

1 **CITE-On: a deep learning approach for online cell identification and trace extraction in**
2 **functional two-photon calcium imaging**

3 Luca Sità^{1,*,#}, Marco Brondi^{1,*,#}, Pedro Lagomarsino de Leon Roig^{1,2}, Sebastiano Curreli¹,
4 Mariangela Panniello¹, Tommaso Fellin^{1,#}

5 ¹ Optical Approaches to Brain Function Laboratory, Istituto Italiano di Tecnologia, Genova, Italy

6 ² University of Genova, Genova, Italy

7

8 * Equal contribution

9

10 # Corresponding authors:

11 Luca Sità, Optical Approaches to Brain Function Laboratory, Istituto Italiano di Tecnologia, Via
12 Morego 30, 16163 Genova, Italy, tel: +39 010 2896549, fax: +39 010 2896230, email: luca.sita@iit.it

13 Marco Brondi, Optical Approaches to Brain Function Laboratory, Istituto Italiano di Tecnologia, Via
14 Morego 30, 16163 Genova, Italy, tel: +39 010 2896549, fax: +39 010 2896230, email:
15 marco.brondi@iit.it

16 Tommaso Fellin, Optical Approaches to Brain Function Laboratory, Istituto Italiano di Tecnologia,
17 Via Morego 30, 16163 Genova, Italy, tel: +39 010 2896549, fax: +39 010 2896230, email:
18 tommaso.fellin@iit.it

19

20

21

22 Keywords: two-photon functional imaging, genetically encoded calcium indicators (GECI), deep
23 neuronal networks, online image segmentation

24

26 **Abstract**

27 *In vivo* two-photon calcium imaging is a powerful approach in experimental neuroscience. However,
28 processing two-photon calcium imaging data is computationally intensive and time-consuming,
29 making online frame-by-frame analysis a challenge. This is especially true for large field-of-view
30 (FOV) imaging. Here, we present CITE-On (Cell Identification and Trace Extraction Online), a
31 convolutional neural network-based algorithm for fast automatic cell identification, segmentation,
32 identity tracking, and trace extraction in two-photon calcium imaging data. CITE-On processes
33 thousands of cells online, including data from mesoscopic two-photon imaging, and provides accurate
34 functional measurements from most neurons in the FOV. Applied to publicly available datasets,
35 CITE-On achieves performance similar to that of state-of-the-art methods for offline analysis.
36 Moreover, CITE-On generalizes across calcium indicators, brain regions, and acquisition parameters
37 in anesthetized and awake head-fixed mice. CITE-On represents a powerful new tool to speed up
38 image analysis and facilitate closed-loop approaches, for example in combined all-optical imaging
39 and manipulation experiments.

40

41 **Introduction**

42 Multi-photon imaging in combination with Genetically Encoded Calcium Indicators (GECI) allows
43 the recording of population activity with high spatial resolution in the intact brain *in vivo*¹⁻⁶. However,
44 multi-photon imaging datasets, in the form of t-series, can be heavy (0.5 GB to > 1 TB) and their
45 processing requires time and computational power. More specifically, the precise identification and
46 segmentation of neuronal structures (typically somata) in a given FOV is critical to extract truthful
47 and reliable information from raw imaging t-series⁶. This step can be complex because of dense
48 GECI staining, low signal-to-noise ratio (SNR), presence of motion artifacts, and large number of
49 neurons in the FOV (e.g., in the case of mesoscopic two-photon imaging^{7,8}).

50

51 Segmentation is typically performed in two ways: *i*) manually, based on visual inspection by an expert
52 user and on selection of pixels into regions of interest (ROIs); *ii*) automatically, employing supervised
53 or unsupervised methods leveraging on static and dynamic properties of the fluorescence signal in
54 the t-series⁹⁻²³. Manual segmentation^{24,25} is time consuming and impractical in case of large datasets
55 and FOVs (e.g. mesoscopic imaging) or when real-time manipulation of experimental conditions is
56 needed^{26, 27}. State-of-the-art automatic approaches apply pixel correlation^{15, 16, 19},
57 principal/independent component analysis (PCA/ICA)^{15, 16}, constrained non-negative matrix
58 factorization (CNMF)^{10, 14, 17}, and deep neural networks (DNN)^{10, 20-22} to perform FOV segmentation.
59 These approaches are usually applied offline and generally take advantage of both neuronal spatial
60 footprints and the temporal dynamics of the fluorescence signal associated to the identified spatial
61 footprints. Consequently, their performance benefits from long acquisitions^{10, 16, 23}, with highly active
62 cells being more easily segmented than rarely active or inactive ones^{10, 22}. Moreover, current methods
63 often require the experimenter to set initialization parameters ahead of the segmentation process^{10,}
64 ^{11, 15-17, 22}. While most of these parameters are generally easy to adjust (*e.g.* frame rate and indicator
65 kinetics), some are inaccessible to the user online (*e.g.* number of expected ROIs in a FOV and spatial
66 constraints on ROI shapes) and must be determined through multiple offline rounds of empirical
67 tuning steps.

68

69 The quality of *in vivo* two-photon calcium imaging is also extremely sensitive to motion artifacts⁶.
70 In particular, the shape and position of imaged cells may change due to motion artifacts correlated
71 with the animal's locomotion, breathing, and heartbeat. In current approaches^{10, 15-17, 22, 28}, successful
72 neuronal segmentation is typically achieved after correcting for motion artifacts: a process requiring
73 additional time and computational power. The output of the segmentation process is thus a static

74 mask, representing the “average” shape and position of each cell throughout the t-series. This
75 approach is impractical whenever cells should be tracked online on a frame-by-frame basis, for
76 instance when a neuronal ensemble (i.e. a group of coactive neurons) must be optogenetically
77 manipulated after being identified^{29, 30}. In fact, neuronal ensembles are dynamic and different cells
78 may belong to a given ensemble at a certain time instant, making it difficult to define *a priori* the
79 neuronal identities belonging to future ensembles³¹. Finally, downstream of segmentation, the
80 dynamic fluorescence signal from each cell must be extracted and “decontaminated” from
81 background signal^{6, 32}. Different approaches are available to this end^{10, 11, 15-17, 32, 32, 33}, all requiring
82 additional computational time. As a result of all these analytical steps, a total processing time of 30
83 to 90 minutes was reported for most efficient methods when processing FOVs of about 500 x 500
84 μm^2 containing hundreds of cells imaged over tens of thousands frames^{10, 11, 22}. Altogether, current
85 analytical approaches are: *i*) still limited in their ability to perform online analysis, which is necessary
86 for closed-loop experiments; *ii*) biased against the identification of rarely active or inactive cells,
87 which could be as informative to target as more active neurons, for example, in longitudinal all-
88 optical imaging and manipulation approaches; *iii*) not validated on large FOVs, such as those
89 generated by mesoscopic imaging.

90

91 Here, we describe CITE-On, a CNN-based algorithm trained to perform neuronal somata
92 identification in two-photon imaging recordings, combined with a light-weight dynamic
93 segmentation and trace extraction pipeline. CITE-On identified hundreds to thousands of neuronal
94 cell bodies on a frame-by-frame basis. Moreover, CITE-On also identified inactive neurons, removing
95 biases towards highly active cells. CITE-On’s light architecture and processing strategy allowed, for
96 the first time, fast automatic segmentation, tracking, and trace extraction in mesoscopic two-photon
97 imaging t-series.

99 **Results**100 *CITE-On: structure and analysis pipeline*

101 CITE-On accepted individual frames from two-photon calcium imaging t-series (Fig. 1a) and it
102 comprised two main parts: an image detector based on the publicly available convolutional neural
103 network (CNN) RetinaNet³⁴ dedicated to the identification of neuronal somata and a custom-built
104 downstream light-weight analysis pipeline, designed for functional trace extraction. The image
105 detector and the analysis pipeline operate as asynchronous parallel processes, in order to provide
106 discrete cell detection update (up to 10 Hz, see text below) and faster than real time functional traces
107 (available at 100 Hz under all experimental conditions tested in this study, see text below). CITE-On
108 required three pre-processing steps ahead of the CNN image detection: *i*) frame downsampling; *ii*)
109 image upscaling; *iii*) replication of the input image into three identical channels (Fig. 1a). The frame
110 downsampling value was set according to the image SNR, while the upscaling factor depended on
111 the ratio between the FOV surface and the average surface of the neuronal somata (see Methods and
112 text below). The length of the frame downsampling window, as well as the value of the upscaling
113 factor, was defined *a priori* and adjusted according to the data to maximize performance. The three
114 identical images were sent to the CNN for image detection (blue rectangle in Fig. 1). The output of
115 the CNN (yellow highlights in Fig. 1) was a set of boxes tightly surrounding each detected cell soma
116 (“bounding boxes”, represented as green squares over the FOVs in Fig. 1a-d). Coordinates and
117 surfaces of each bounding box were used in the analysis pipeline, to generate: *i*) cell identity
118 assignment and tracking along the t-series; *ii*) dynamic segmentation of neuronal somata; *iii*)
119 background subtraction procedure and extraction of neuronal functional trace.

121 CITE-On worked both offline, after the acquisition was completed and the whole t-series was
122 available (Fig. 1b), or online, using individual imaging frames as inputs, continuously streamed from
123 the experimental set-up during the acquisition of the imaging t-series (Fig. 1c-d). In the offline
124 pipeline, a Fourier-transform approach³⁵ was used to correct for planar motion artifacts throughout
125 the t-series (Fig. 1b). The frame downsampling corresponded to the projection of the entire t-series
126 onto its temporal median calculated across all frames (Fig. 1b). Soma detection was then performed
127 once on the pre-processed median image (detection in Fig. 1b), and bounding boxes were generated
128 for each frame of the t-series (Fig. 1b, yellow highlights). Each bounding box was associated with a
129 score, representing network confidence in cell detection. Bounding boxes with intersection over union
130 (IoU) < 20 % were considered as separate neuronal identities. When IoU of two bounding boxes was
131 > 20 %, the bounding box with the highest score was retained.

132

133 In the online pipeline, no motion correction was performed, and the user selected between two
134 downsampling strategies depending on the SNR of the data, and on the required upscaling factor in
135 the pre-processing step. In case of relatively high SNR and low upscaling factors (Fig. 1c), a sliding
136 average was calculated on the first n frames of the t-series and updated with every new individual
137 frame generated by the microscope. Neuronal detections were updated for each imaging frame
138 starting from the $n+1^{\text{th}}$ frame. When the SNR was relatively low and the upscaling factor large (Fig.
139 1d), a step average approach was performed, where the input for the image detector was the average
140 projection of blocks of n frames. Additional n frames were thus required for generating the next step
141 average projection and the detections were updated every n frames. The maximum detection rate
142 decreased with the upscaling factor, with a peak rate of 10 Hz with upscaling factor equal to 1 (Fig.
143 1e). Active detections (*i.e.* detections in the current sliding average or step average) and past
144 detections (*i.e.* detections in any previous sliding or step average) were continuously tracked and

145 updated (Fig. 1c-d, detection update and I.D. tracking). Specifically, active detections were compared
146 with past detections at each step of detection update and a new identity was added (and included in
147 the tracking system) every time the surface of an actively detected bounding box had $\text{IoU} < 25\%$
148 with any of the previously identified boxes. Bounding boxes from active detections with $\text{IoU} > 25\%$
149 with those of past detections did not change identities of previously detected boxes, but their positions
150 and shapes were updated according to the position in the most recent detection step. All past
151 detections without updates were retained in the tracking system in the form of their last active
152 detection for the remaining part of the t-series (Supplementary Movie 1).

153

154 For both the online and offline pipelines (Fig. 1f-g), bounding boxes were used to generate a dynamic
155 segmentation of the t-series and to identify ROIs. The distribution of fluorescence values inside each
156 bounding box was computed at each frame (Fig. 1f, left). Only pixels with values between the 80th
157 and the 95th percentile of the box's fluorescence distribution were assigned to the ROI corresponding
158 to the cell soma (white pixels of the binary mask in Fig. 1f, right). Since pixel assignment to cell
159 somata in each individual box was updated at each frame, the resulting dynamic segmentation was
160 updated online for every new frame. All the FOV pixels that were not included in any bounding box
161 were assigned to a global background ROI. The fluorescence intensity of all pixels belonging to the
162 global background ROI was averaged at each frame to obtain the background signal (*bg*). Moreover,
163 at each frame, the *bg* was subtracted from the fluorescence of each segmented neuronal ROI,
164 generating functional fluorescence traces (Fig. 1g). Since shape, number, and position of bounding
165 boxes changed as the t-series progressed (according to active detections and tracking), the pixels
166 assigned to *bg* also changed in number and identity across frames. Identity tracking, segmentation,
167 and functional trace extraction required on average 10 ms *per* frame (either offline or online).

168

169 *Ground truth generation and training of the image detector*

170 The ResNet50 Feature Extractor CNN incorporated in CITE-On was not originally developed for
171 detecting neuronal somata, rather for the analysis of natural images and it was trained on > 1 million
172 RGB images across 80 classes (<http://www.image-net.org/>). We decided to use a transfer learning
173 strategy³⁶ to adapt this efficient detection architecture to the identification of neuronal somata (*i.e.* a
174 single class) in grey-scale two-photon images. This choice was dictated by the fact that available two-
175 photon calcium imaging datasets (<http://neurofinder.codeneuro.org/>, <http://help.brain->
176 map.org/display/observatory/Data++Visual+Coding) are far too small for an *ab-initio* CNN training.
177 Moreover, they are too homogeneous in terms of calcium indicator used, FOV dimensions, cell
178 density, acquisition frame period, SNR and background signal contamination²², making them
179 suboptimal even for a transfer learning strategy. For example, no publicly available large dataset
180 comprises imaging data collected using red-shifted GECIs, such as jRCaMP1a. We thus decided to
181 use a dedicated dataset for training and internal validation. In this way, we employed publicly
182 available datasets to test CITE-On performance and its generalization capability on never-before-seen
183 data. The dedicated dataset included 197 t-series from 28 mice acquired using different acquisition
184 parameters (see Methods). More specifically, we included 121 t-series from layer IV neurons of the
185 somatosensory cortex expressing either GCaMP6f, GCaMP6s, or GCaMP7f (globally indicated as
186 “LIV”) and 76 t-series from the CA1 pyramidal neurons of the hippocampus expressing both
187 GCaMP6f and jRCaMP1a (indicated as “CA1 GCaMP6f” and “CA1 jRCaMP1a”, respectively). We
188 included t-series with heterogeneous median fluorescence and SNR in order to reduce potential biases
189 toward bright cells during the training process, while avoiding large differences between groups of
190 data that could have generated better performance on specific subsets of t-series (Supplementary Fig.
191 1a, b).

192

193 To obtain a consensus Ground Truth (GT) annotation of the t-series used for training and validation
194 of the CNN, two human graders manually annotated all t-series, defining the tightest rectangular box
195 fitting each visible cell in each FOV. Manual GT annotation was preferred to automatic segmentation
196 for two main reasons: *i*) available automatic segmentation approaches rely on both functional (i.e.
197 fluorescence signal dynamics across frames) and morphological features⁹⁻²³, while we wanted the
198 GT annotation to be exclusively based on morphological features (see below); *ii*) manual annotation
199 is still frequently considered more accurate than automatic methods^{10, 22}. Initially, manual annotation
200 on single frames by two graders produced only few neuronal identities because cells were only visible
201 in a minority of frames (Supplementary Fig. 1c, Supplementary Movie 2). This could be due to the
202 low basal emission of some of the indicators used (e.g. GCaMP6f), to the variable expression level
203 of the calcium indicator across cells, and to the sparse activity profile of the imaged cells. In order to
204 increase the visibility of neurons, we created high contrast single images representative of each t-
205 series. To this end, we first corrected each t-series for planar motion artifacts, and then collapsed each
206 acquisition onto its median projection (Supplementary Fig. 1d). These images were sharpened
207 (Supplementary Fig. 1e) and gamma corrected. Brightness and contrast were adjusted in order to
208 obtain a distribution of intensity values spanning the whole bit range. The sharpened images, named
209 enhanced median projection (EMP) (Supplementary Fig. 1f, g), were used for manual annotation
210 (Supplementary Fig. 1h, Supplementary Movie 3). In training and validation datasets, grader #1
211 annotated 14,425 boxes, while grader 2 annotated 12,912 (Supplementary Table 1). The bounding
212 boxes produced by grader #1 and grader #2 and their superposition in different experimental
213 preparations are shown in Supplementary Fig. 2. Annotations from the two graders largely overlapped
214 (weighted average precision or mAP, 0.77 ± 0.08 ; F-1 score, 0.93 ± 0.02 ; precision: 0.98 ± 0.01 ;
215 recall: 0.88 ± 0.12 , N = 197 EMPs, Supplementary Table 2). Given the high similarity of the
216 independent annotations provided by the two graders, we defined the consensus GT for our entire

217 dataset as the union of the two sets of annotations. Given that our dataset contained partially
218 overlapping FOVs and more than one t-series acquired on the same FOV, the dataset was manually
219 split into training (160 t-series) and validation (37 t-series) subsets. To avoid data leakage and to
220 decrease overfitting³⁷, t-series from the same or similar FOVs were included in either the training
221 dataset or the validation dataset. We trained the CITE-On image detector on the training dataset and
222 evaluated its performance on the validation dataset.

223

224 *Performance of the image detector*

225 We trained the CITE-On image detector on our consensus GT annotations achieving the best
226 performance after 17 training epochs (mAP: 0.79 on the validation dataset). We first evaluated CITE-
227 On performance using the offline pipeline on the validation dataset. A representative CITE-On output
228 for a CA1-GCaMP6f, a CA1-jRCaMP1a, and a LIV t-series is shown in Fig. 2a-d and Fig. 2e,
229 respectively. For the whole validation dataset, the mean \pm s.d. F-1 score, Precision, and Recall are
230 reported in Fig. 2f and Supplementary Table 3.

231

232 We then used the online pipeline and calculated the F-1 of CITE-On detections on the motion
233 corrected validation dataset, which was used as input to CITE-On at the actual frame rate occurring
234 during acquisition. Our validation data required an upscaling factor of 1, compatible with a maximum
235 detection rate of 10 Hz, while acquisition frame rates varied between 1.5 Hz for LIV and 3 Hz for
236 CA1 acquisitions. We empirically explored the effect of frame downsampling on detection
237 performance using the sliding average approach with different numbers of frames (n). Our aim was
238 to maximize F-1 and score threshold for detections, while minimizing n . F-1 increased with n between
239 1 and 20 frames, value at which the absolute maximum was observed (Fig. 3a). Using the sliding
240 average, an initial delay of 6.6 s for LIV data and 14 s for CA1 data was necessary before CITE-On

241 processed each frame in real time at a detection rate of 10 Hz. F-1 values calculated on sliding
242 averages of 20 frames are reported in Fig. 3a. Detections (green boxes) are displayed together with
243 GT annotations (red boxes) for both the GCaMP6s (Fig. 3b, left) and the jRCaMP1a channels (Fig.
244 3b, middle). The superimposition of the detections from both channels is shown on Fig. 3b (right).
245 Similarly, in Fig. 3c, we show GT and online detections (red and green boxes, respectively) for a
246 representative LIV (GCaMP6s) t-series.

247

248 To quantify the impact of motion artifacts on online detection accuracy, we calculated the F-1 score
249 on validation t-series that were not corrected for motion artifacts (Supplementary Table 4). To this
250 end, we translated the GT annotations for each frame according to the shift vectors produced by the
251 motion correction algorithm implemented when building the relative EMPs (see Methods). We
252 observed no significant difference between the F-1 obtained on motion corrected and non-corrected
253 validation t-series (Fig. 3d), suggesting that the motion correction step was not necessary to achieve
254 higher performance with our approach. The F-1 values across frames for the non-motion-corrected
255 data are reported in Fig. 3e. The distribution of motion displacements is shown in Supplementary Fig.
256 3.

257

258 *Online data processing*

259 We developed a fast method to dynamically segment each detected cell based on the corresponding
260 bounding box, relying on the instantaneous (i.e. frame-wise) fluorescence statistics of the pixels
261 inside each box (Fig. 4, Supplementary Movie 4). Specifically, the fluorescence intensity distribution
262 of the pixels inside each bounding box was first computed frame-wise. Pixels with fluorescence
263 values between the 80th and 95th percentile of the distribution were then selected as belonging to
264 neuronal somata. The values of selected pixels were averaged, and the resulting fluorescence trace

265 was “denoised” by subtracting, at each frame, the *bg* signal. This simple method was computationally
266 light, an important requirement to achieve fast frame-by-frame data processing (trace extraction rate,
267 100 Hz). Bounding boxes detected by CITE-On on a representative LIV t-series and a representative
268 CA1 t-series are shown in Fig. 4a. Representative fluorescence traces extracted by CITE-On on the
269 two t-series are displayed in Fig. 4b-c. Fig. 4d-e shows the cross correlation matrix (lower left
270 triangle) and the dendrogram analysis (upper right triangle) of all the identified cells before (Fig. 4d-
271 e left) and after (Fig. 4d-e, right) *bg* subtraction.

272

273 We compared traces extracted by CITE-On with those extracted by CaImAn, a state-of-the-art
274 method based on CNMF¹⁰. We used the bounding boxes generated offline by CITE-On to build
275 binary masks that were used as seeds to initialize the seeded-CNMF algorithm¹⁰. The spatial
276 components of the CNMF were non-zero only inside the bounding boxes identified by CITE-On,
277 therefore, the detected factors from seeded-CNMF had one-to-one correspondence with the detected
278 boxes from CITE-On. We first observed very high pairwise cross correlations between the *bg* traces
279 extracted with the two methods (Supplementary Fig. 4a). We then asked how the *bg* signal calculated
280 over the whole FOV (*bg*) correlated with the ‘local *bg*’, that is, the background activity calculated for
281 each cell from the pixels in the immediate surroundings of the relative bounding box (see Methods).
282 The average correlation between *bg* and local *bg* traces was high (Supplementary Fig. 4b). Given the
283 high correlation values observed between *bg* and local *bg* and the lower computational cost of *bg*, we
284 decided to implement the *bg* method only. We then tested how the *bg*-subtracted functional traces
285 calculated by CITE-On compared with those extracted with seeded-CNMF and we again observed
286 high correlation values (Supplementary Fig. 4c).

287

288 Although the average correlation values for the cells extracted with the two methods were high, some
289 neuronal pairs showed lower correlations. We asked whether the low correlation values emerged from
290 pairs of cells with low SNR. We found that pairwise correlation values for traces extracted with the
291 two methods increased with the SNR of the corresponding neuronal traces (Supplementary Fig. 4f),
292 indicating that indeed the functional traces obtained with the two methods were more similar when
293 the trace SNR was high.

294

295 *Offline performance on never-before-seen recordings*

296 To test the robustness of our image detection approach and its ability to generalize across
297 experimental conditions, we tested CITE-On on three additional datasets, which were not used during
298 the training and validation phases. The three datasets were: the Allen Brain Observatory repository
299 (ABO, 19 t-series divided into 9 superficial, ABO_{sup}, t-series, acquired in visual cortex at depth 175
300 μm, and 10 deep, ABO_{deep}, t-series, acquired in visual cortex at depth 275 μm), the Neurofinder
301 Challenge dataset (28 t-series, divided into 19 t-series, NF_{train}, 9 t-series, NF_{test}, from different
302 preparations at different depths), and a dataset acquired in our laboratory using GRIN-based
303 endoscopic two-photon imaging of the ventral posteromedial thalamic nucleus (VPM, 9 t-series). The
304 three datasets were first manually annotated *de novo* to obtain the GT annotation (Supplementary Fig.
305 5, Supplementary Table 1-2). Because the ratios between FOV and cell surface were variable across
306 the ABO, NF, and VPM datasets, and different from our validation dataset, we optimized the
307 upscaling factor and used the one that maximized the F-1 score (Fig. 5a-e). The offline performance
308 (Precision, Recall and F-1 score) obtained using optimized upscaling factors for each dataset is shown
309 in Fig. 5f and Supplementary Table 3.

310

311 We compared the offline detection performance of CITE-On (Supplementary Table 3) with state-of-
312 the-art alternative segmentation approaches such as STNeuroNET ²², CaImAn_{online}¹⁰, CaImAn_{batch}¹⁰,
313 Suite2P ¹⁶, HNCcorr ¹⁹, UNet2DS ²³ on the ABO and NF datasets provided in ²². We used our GT
314 annotations to test CITE-On and those reported in ²² to test the alternative methods. We found that
315 the average CITE-On performance was equivalent to that computed with other methods (Fig. 6).
316 Moreover, the background signals calculated using seeded CNMF and CITE-On presented high cross
317 correlation values for all three datasets (Supplementary fig. 6a). Similarly, high cross correlation was
318 measured between local *bg* and *bg* signals (Supplementary Fig. 6b), as well as between CITE-On
319 extracted functional traces after global *bg* subtraction and CaImAn extracted traces (Supplementary
320 Fig. 6c). As for our previous characterization on the validation dataset, cross correlation values of the
321 extracted traces with CITE-On and CaImAn increased with SNR in all datasets (Supplementary Fig.
322 6f).

323

324 We compared true positive detections obtained with CITE-On with ABO true positive detections and
325 with STNeuroNET true positive detections (Fig. 7a-c). Examples of cells identified by CITE-On and
326 present in the true positive detections from both ABO and STNeuroNET are reported in Fig. 7 d,
327 together with their relative background-subtracted traces extracted by CITE-On. We analyzed the
328 functional traces of all the CITE-On true positives (i.e. including cells detected in ABO and
329 STNeuroNET) after background subtraction (Fig. 7f). Similarly to what described for our validation
330 dataset, the dendrogram sorting showed blocks with different cross correlation values for various
331 subgroups of cells. On average, the number of CITE-On detected cells exceeded the number of
332 identities available in public repositories (Supplementary Fig. 7, Supplementary Table 1).

333

334 Since state-of-the-art segmentation methods have a bias against inactive cells, we investigated
335 whether CITE-On-only cells were inactive. We quantified the number of calcium events *per* detection
336 in the ABO true positives (ABO TP), STNeuroNET true positives (STNuroNET TP), CITE-On true
337 positives (CITE-On TP), and CITE-On-only true positives (CITE-On-only, Supplementary Fig. 7).
338 The number of cells with few detected calcium events was larger for CITE-On TP (Supplementary
339 Fig. 7). Moreover, although some of the CITE-On-only cells were silent (as expected), a large fraction
340 of them displayed detectable activity (in the whole ABO dataset, 91 % displayed at least one calcium
341 event and 69 % showed at least ten calcium events, Supplementary Fig. 7a-d). Moreover, the structure
342 of the dendrogram built with CITE-On-only cells resulted similar to the one obtained with the
343 complete set of CITE-On true positives (Fig. 7g), indicating that the group of active cells detected
344 exclusively by CITE-On recapitulated the functional structure of most active neurons. The number
345 of detected calcium events with CITE-On (all detections, true positives, false positives, and CITE-
346 On-only identities), ABO (true and false positives), and with STNeuroNET (true and false positives)
347 is reported in Supplementary Fig. 7e. The number of detections in the different datasets is reported in
348 Supplementary Table 5.

349

350 *Online performance on never-before-seen recordings*

351 We ran CITE-On online using our GT annotation on each frame of the ABO, NF, and VPM datasets.
352 When computing the F-1 score, we tested different sizes of the sliding average, in order to define the
353 smallest number of frames required to achieve real-time processing. The absolute maximum F-1 score
354 was achieved in 10 frames for the ABO dataset (both ABO_{sup}, and ABO_{deep}), 20 frames for the VPM
355 dataset, and 200 frames for the NF (train and test) dataset (Fig. 8a-e). For the ABO dataset, 0.3 s
356 were necessary to acquire 10 frames (0.00086 % of the whole time series of average duration 115,635
357 \pm 130 frames) with a CITE-On detection rate of 5 Hz and an upscaling factor of 2. The NF dataset

358 required upscaling factors of 2.4 and 2.6 for NF_{train} and NF_{test} , respectively. A time window of 28.5 s
359 was required to acquire the 200 frames necessary to reach peak F-1 value (5.4 % of the whole t-series
360 of average duration $3,697 \pm 1,874$ frames) with a CITE-On detection rate of 4 Hz for NF_{train} and 3 Hz
361 for NF_{test} . Given the relatively low online performance for these latter datasets, we decided to measure
362 the F-1 score using an alternative frame downsampling strategy: step average. Using this approach,
363 we found that detection performance remained high (Fig. 8i and Supplementary Table 4). The F-1
364 score calculated online as a function of the length of the t-series for ABO (ABO_{sup} and ABO_{deep}
365 together), NF (NF_{train} and NF_{test}) and VPM data is shown in Fig. 8j. Stable F-1 scores were observed
366 within 30 % of the total length of the t-series. Processing time required for running the online pipeline
367 on the ABO, VPM, and NF datasets was not different from what described previously for our
368 validation datasets.

369

370 *Analysis of large FOV mesoscopic images*

371 Given the speed and light-weight of the CITE-On architecture, we tested if it could be applied for
372 detecting cells in the mesoscopic imaging t-series described in⁷. Because the dimensions of the input
373 image were too large (1792 pixels x 1682 pixels, or 4.8 mm x 4.8 mm, 1 μm pixel size) to fit the CNN
374 architecture using the appropriate upscaling factor based on the FOV/neuron surface ratio (upscale
375 factor for direct processing, 12.8), we tiled the entire mesoscopic FOV in 272 subfields (subfield
376 dimension, 128 pixels x 128 pixels, 28 pixels overlap). Each subfield was appropriately upscaled
377 (upscale factor, 1; score threshold, 0.4; detection rate, 10 Hz). To increase speed, we multiplexed
378 the CNN detector process and processed subfields in batches of 8 images in parallel, until completion.
379 Identity duplicates were suppressed using a non-maximum suppression algorithm, where boxes
380 having an IoU > 25% were considered duplicates and only the one with the highest score was retained
381 (see Methods). Image detection outputs were finally recombined to reconstruct the entire FOV. In

382 Fig. 9a, we report CITE-On detected bounding boxes on the entire FOV obtained using the offline
383 pipeline. Two FOV patches are magnified (Fig. 9b) to show the shape of the detected cells (total
384 number of detected somata, 4,842). In Fig. 9c, we show representative fluorescence traces obtained
385 with CITE-On from five cells from the whole mesoscopic FOV. Fig. 9d shows the cross correlation
386 matrix (lower left triangle) of all the identified cells for the first 700 frames of the t-series. The
387 dendrogram analysis (upper right triangle of Fig. 9d) highlighted several distinct functional modules
388 observed in the identified neuronal population.

389

390 The single subfields were processed by the CITE-On detector at 10 Hz (upscale factor, 1), while
391 segmentation, tracking, and functional trace extraction were performed at 100 Hz (see previous
392 results). Parallel processing of all the 272 sub-fields required 12.6 s for each detection step. Therefore,
393 with a step average downsampling approach including 25 frames (13.2 s of mesoscopic imaging time
394 since acquisition rate was 1.9 Hz), and while extracting traces faster than the incoming frames, we
395 minimized the CITE-On detection lag with respect to the running acquisition. With this strategy, we
396 achieved an online F-1 score of 0.54 with a score threshold of 0.25 (quantified on four patches from
397 the entire FOV). CITE-On can thus be efficiently applied for fast processing of mesoscopic two-
398 photon t-series.

399

400 **Discussion**

401 In this study, we developed a CNN-based algorithm, CITE-On, for fast analysis of two-photon
402 imaging recordings. CITE-On performed online accurate identification of neuronal somata, tracking
403 of identities across frames, dynamic segmentation, and functional trace extraction with background
404 subtraction. CITE-On generalized across calcium indicators, brain regions, acquisition parameters,
405 and it was successfully applied to data obtained using different surgical and optical preparations (e.g.

406 chronic superficial imaging window, chronic deep imaging window, and endoscopic GRIN lens-
407 based deep imaging).

408

409 Our image detection strategy was based on RetinaNET, a CNN originally developed to detect natural
410 images ³⁴. On one side, this choice was justified by RetinaNET's excellent performance in object
411 recognition and by its readiness. On the other hand, it required us to exploit RetinaNET on a set of
412 images, grey-scale two-photon fluorescence images of neurons, that were remarkably different from
413 those RetinaNET was originally trained on. To compensate for this difference, and to have a large
414 and heterogeneous dataset for training and validation of the detection algorithm, we built a dedicated
415 library of hundreds of two-photon imaging t-series acquired with different GECIs, in different regions
416 of the mouse brain, at different frame rates, using different surgical/optical preparations, and showing
417 variable image quality. In this dataset, a reliable GT consensus was reached using the annotations of
418 two human graders, which allowed us to evaluate CITE-On performance. To obtain a GT annotation
419 insensitive to the graders' biases, a potential alternative approach could have been to generate an *in*
420 *silico* dataset for network training and validation ³⁸. No online libraries of this kind are currently
421 available, but we foresee that this approach may represent an extremely useful method to optimize
422 future CNN-based approaches for the analysis of two-photon functional data. It is also worth noticing
423 that utilizing public datasets already used for training and validation of alternative processing
424 toolboxes ²² would have given us the possibility to take advantage of third-party GT annotations.
425 However, we decided not to do so because: *i*) CITE-On would have likely inherited the annotation
426 bias toward more active cells, which is shared by existing publicly available repositories; *ii*) by using
427 public datasets exclusively for the validation of the CITE-On image detector (rather than for training,
428 too), we avoided any chance of data leakage, and we demonstrated that CITE-On generalized well to
429 never-before-seen data.

430

431 CITE-On performed as state-of-the-art algorithms^{10, 11, 16, 19, 22, 39} on publicly available datasets and,
432 importantly, it did so in a much shorter time. In fact, only a few seconds were needed to have online,
433 frame-by-frame, accurate ROI segmentation, identity tracking, *bg* subtraction, and functional trace
434 extraction. Four main characteristics were crucial for CITE-On's high performance. First, CITE-On
435 relied exclusively on morphological features to identify neurons. Second, neuronal identification was
436 dynamic and it adapted to changes in shape, position, and activity of the detected cells frame-by-
437 frame, avoiding time-expensive motion correction procedures. Third, once bounding boxes were
438 identified in individual frames, we used a simple cost-effective strategy to extract pixels belonging to
439 neuronal ROIs based on pixel's brightness. Fourth, we implemented a fast background subtraction
440 strategy, limiting computational costs. When applied in the online modality, these characteristics were
441 crucial to achieve real-time frame-by-frame trace extraction, something current approaches do not
442 achieve⁹⁻²⁵, while maintaining high cell detection performance.

443

444 The most common automatic analytic pipelines for the analysis of two-photon imaging t-series
445 typically rely on both spatial features and fluorescence signal fluctuations⁹⁻²³. This necessarily
446 requires the use of several, if not all, frames in a t-series, precluding efficient online analysis. In
447 contrast, CITE-On was conceived such that the analysis pipeline was performed frame-by-frame
448 based on morphological features. This dynamic analytical process was continuously updated over
449 time allowing the algorithm to cope with planar motion artifacts, a feature that, unlike other available
450 methods, enabled us to bypass the motion correction pre-processing step and to save further
451 computational time. Importantly, neuronal identities identified frame-by-frame were reliably tracked
452 over the t-series. To extract pixels belonging to neuronal ROIs within a bounding box, we used a
453 criterion based on pixel brightness and to correct for background contamination we devised a global

454 rather than local neuropil computation. These procedures had the advantage of being simply and
455 simultaneously implementable for thousands of cells, allowing the extraction of functional traces at
456 low computation cost. The observation that functional fluorescence traces extracted by CITE-On were
457 highly correlated with those extracted on the same bounding boxes by a state-of-the-art method, i.e.
458 CaImAn¹⁰, confirmed the validity of our computationally effective approach.

459

460 Thanks to the features described above, CITE-On efficiently processed full mesoscopic two-photon
461 t-series (FOV dimension, 5 mm x 5 mm). CITE-On analyzed full mesoscopic images dividing each
462 image in subfields and processing subfields in parallel. CITE-On's detector processed single subfields
463 at 10 Hz, while segmentation, tracking, and functional trace extraction were performed at 100 Hz.
464 Parallel processing of all the 272 sub-fields generating a whole mesoscopic FOV required 12.6 s for
465 each detection step. Thus after 12.6 s, trace extraction could be performed at 100 Hz on thousands
466 cells. Besides its application online, the offline application of CITE-On is also going to be extremely
467 powerful for the identification of the thousands neurons imaged in mesoscopic two-photon functional
468 imaging.

469

470 Closed-loop all-optical experiments are fundamental to investigate whether models of network
471 dynamics, circuit connectivity, and causality are accurate²⁶. Recently, all-optical closed-loop
472 experiments have been validated²⁷. For example, using this approach specific groups of neurons were
473 activated based on the readout of ongoing activity in a reference cell. However, the closed-loop
474 strategy described in²⁷ was based on *a priori* identification of the reference cell. Because CITE-On
475 allows efficient frame-by-frame cell identification and trace extraction, it will enable a new type of
476 experiment in which the loop is closed based on real time identification and readout of any neuron or
477 group of neurons in the FOV.

478

479 CITE-On differs in several main features from OnACID¹¹, a recent method that combines CNMF
480 and fast deconvolution for the online analysis of calcium data. First, CITE-On does not need any
481 initialization step in the online pipeline. In contrast, OnACID requires the user to provide the expected
482 number of cells to be identified, and it starts the online analysis based on the CNMF output of an
483 initial portion of the t-series (typically 1000-3000 frames), generating an initial temporal lag of 2-4
484 minutes¹¹. Second, CITE-On does not require the correction of motion artifacts in the online pipeline,
485 saving further computational time. Third, CITE-On performs tracking, dynamic segmentation, and
486 functional trace extraction processes at 100 Hz independently of the number of detected neurons and
487 their activity. This feature allows maintaining high online performance on FOVs characterized by
488 large number of neurons (e.g., those obtained with mesoscopic two-photon imaging) and sparse
489 activity. Four, CITE-On does not use local pixel correlation for cell identification, which may be
490 advantageous to separate nearby synchronous cells. In contrast, the OnACID fast deconvolution
491 approach may be more efficient in separating adjacent cells with different temporal profiles of
492 fluorescence emission¹¹. Finally, OnACID was tested on two datasets with rather homogenous
493 acquisition parameters¹¹. Here, we demonstrate that CITE-On generalizes across indicators (e.g.,
494 GCaMP6s, GCaMP6f, GCaMP7f, and jRCaMP1a) and across data acquired with different pixel size,
495 SNR, and frame rate. This property adds flexibility in the application of online analysis tools to
496 different experimental conditions.

497

498 Because cell identification was based only on spatial features, CITE-On identified both active and
499 silent cells. This unique characteristic of CITE-On is important because it adds further flexibility in
500 designing imaging experiments. Neurons that are silent in a t-series may change their level of activity
501 in subsequent acquisitions depending on the behavioral state of the animal or because of external

502 manipulations⁴⁰⁻⁴². Thus, being able to track cells regardless of their activity level is key, for instance,
503 for investigating the sensory information beard by neurons that dramatically change their activity
504 throughout longitudinal imaging experiments. Biasing the analysis toward active neurons, as
505 currently done by most approaches, intrinsically skews the proportion of cells that are responsive to
506 a given stimulation in a certain brain region. In this regard, it is also interesting to note that neurons
507 that were detected only by CITE-On and not by other state-of-the-art approaches comprised silent
508 cells but, unexpectedly, also active cells. Active CITE-On only cells may be associated with low SNR
509 of fluorescence signals and they may be missed by current approaches¹⁰.

510

511 In summary, we developed CITE-On, a new tool to effectively process two-photon imaging data
512 frame-by-frame, while maintaining similar cell detection and trace extraction performance of existing
513 offline state-of-the-art methods. Future developments of CITE-On will likely include its optimization
514 for one-photon imaging^{43, 44}, its application to genetically encoded voltage indicators⁴⁵ as well as to
515 volumetric two-photon imaging⁴⁶.

516

517 **Author contributions**

518 LS, MB, PL developed software and performed analysis. MB, SC performed experiments. LS, MB,
519 TF conceived the project. TF coordinated the project. LS, MB, MP, TF wrote the paper with inputs
520 from other authors. All authors approved the final version of the manuscript.

521

522 **Acknowledgments**

523 We thank Dr. C. Arlt and Dr. C. Harvey for sharing mesoscopic imaging data, Dr. A. Sattin for sharing
524 VPM recordings, Dr. S. Fiorini and Dr. A. Barla for discussion and comments on algorithm
525 development, Dr. D.S. Kim and the GENIE project and Dr. J.M. Wilson for constructs, Dr. S. Succol

526 for technical support, and Dr. D. Vecchia for help with the figures. This work was supported by the
527 European Research Council (ERC, NEURO-PATTERNS) and NIH Brain Initiative (U19
528 NS107464).

529

530 **Declaration of interest**

531 The authors declare no competing interests.

532

533 **References**

- 534 1. Helmchen,F. & Denk,W. Deep tissue two-photon microscopy. *Nat Methods* 2, 932-
535 940 (2005).
- 536 2. Wang,T. *et al.* Three-photon imaging of mouse brain structure and function through
537 the intact skull. *Nat. Methods* 15, 789-792 (2018).
- 538 3. Chen,T.W. *et al.* Ultrasensitive fluorescent proteins for imaging neuronal activity.
539 *Nature* 499, 295-300 (2013).
- 540 4. Dana,H. *et al.* Sensitive red protein calcium indicators for imaging neural activity.
541 *Elife.* 5, 12727 (2016).
- 542 5. Svoboda,K. & Yasuda,R. Principles of two-photon excitation microscopy and its
543 applications to neuroscience. *Neuron* 50, 823-839 (2006).
- 544 6. Harris,K.D., Quiroga,R.Q., Freeman,J., & Smith,S.L. Improving data quality in
545 neuronal population recordings. *Nat Neurosci.* 19, 1165-1174 (2016).
- 546 7. Sofroniew,N.J., Flickinger,D., King,J., & Svoboda,K. A large field of view two-
547 photon mesoscope with subcellular resolution for in vivo imaging. *Elife.* 5, (2016).
- 548 8. Tsai,P.S. *et al.* Ultra-large field-of-view two-photon microscopy. *Opt. Express* 23,
549 13833-13847 (2015).

- 550 9. Diego,F. & Hamprecht,F. Sparse space-time deconvolution for Calcium image
551 analysis. *Prooceedings of the 27th International Conference on Neural Information
552 Processing Systems* 1, 64-72 (2014).
- 553 10. Giovannucci,A. *et al.* CaImAn an open source tool for scalable calcium imaging data
554 analysis. *Elife*. 8, (2019).
- 555 11. Giovannucci,A. *et al.* OnACID: Online analysis of calcium imaging data in real time.
556 *Advances in Neural Information Processing Systems*(2017).
- 557 12. Guan,J. *et al.* NeuroSeg: automated cell detection and segmentation for in vivo two-
558 photon Ca(2+) imaging data. *Brain Struct. Funct.* 223, 519-533 (2018).
- 559 13. Kaifosh,P., Zaremba,J.D., Danielson,N.B., & Losonczy,A. SIMA: Python software for
560 analysis of dynamic fluorescence imaging data. *Front Neuroinform.* 8, 80 (2014).
- 561 14. Maruyama,R. *et al.* Detecting cells using non-negative matrix factorization on calcium
562 imaging data. *Neural Netw.* 55, 11-19 (2014).
- 563 15. Mukamel,E.A., Nimmerjahn,A., & Schnitzer,M.J. Automated analysis of cellular
564 signals from large-scale calcium imaging data. *Neuron* 63, 747-760 (2009).
- 565 16. Pacitariu,M. *et al.* Suite2p: beyond 10,000 neurons with standard two-photon
566 microscopy. *BioRxiv*(2017).
- 567 17. Pnevmatikakis,E.A. *et al.* Simultaneous Denoising, Deconvolution, and Demixing of
568 Calcium Imaging Data. *Neuron* 89, 285-299 (2016).
- 569 18. Reynolds,S. *et al.* ABLE: An Activity-Based Level Set Segmentation Algorithm for
570 Two-Photon Calcium Imaging Data. *eNeuro*. 4, (2017).
- 571 19. Spaen,Q. *et al.* HNCcorr: A Novel Combinatorial Approach for Cell Identification in
572 Calcium-Imaging Movies. *eNeuro*. 6, (2019).

- 573 20. Stringer,C., Wang,T., Michaelos,M., & Pachitariu,M. Cellpose: a generalist algorithm
574 for cellular segmentation. *Nat. Methods*(2020).
- 575 21. Apthorpe,N.J. *et al.* Automatic Neuron Detection in Calcium Imaging Data Using
576 Convolutional Networks. *arXiv*(2016).
- 577 22. Soltanian-Zadeh,S., Sahingur,K., Blau,S., Gong,Y., & Farsiu,S. Fast and robust active
578 neuron segmentation in two-photon calcium imaging using spatiotemporal deep learning.
579 *Proc. Natl. Acad. Sci. U. S. A* 116, 8554-8563 (2019).
- 580 23. Ronneberger,O., Fischer,P., & Brox,T. U-net:Convolution networks for Biomedical
581 Image Segmentation. *arXiv*(2015).
- 582 24. Ohki,K., Chung,S., Ch'ng,Y.H., Kara,P., & Reid,R.C. Functional imaging with cellular
583 resolution reveals precise micro-architecture in visual cortex. *Nature* 433, 597-603 (2005).
- 584 25. Smith,S.L. & Häusser,M. Parallel processing of visual space by neighboring neurons
585 in mouse visual cortex. *Nat. Neurosci.* 13, 1144-1149 (2010).
- 586 26. Grosenick,L., Marshel,J.H., & Deisseroth,K. Closed-loop and activity-guided
587 optogenetic control. *Neuron* 86, 106-139 (2015).
- 588 27. Zhang,Z., Russell,L.E., Packer,A.M., Gauld,O.M., & Häusser,M. Closed-loop all-
589 optical interrogation of neural circuits in vivo. *Nat. Methods* 15, 1037-1040 (2018).
- 590 28. Dombeck,D.A., Khabbaz,A.N., Collman,F., Adelman,T.L., & Tank,D.W. Imaging
591 large-scale neural activity with cellular resolution in awake, mobile mice. *Neuron* 56, 43-57
592 (2007).
- 593 29. Bovetti,S. & Fellin,T. Optical dissection of brain circuits with patterned illumination
594 through the phase modulation of light. *J. Neurosci. Methods* 241, 66-77 (2015).

- 595 30. Panzeri,S., Harvey,C.D., Piasini,E., Latham,P.E., & Fellin,T. Cracking the Neural
596 Code for Sensory Perception by Combining Statistics, Intervention, and Behavior. *Neuron*
597 93, 491-507 (2017).
- 598 31. Carrillo-Reid,L., Yang,W., Kang Miller,J.E., Peterka,D.S., & Yuste,R. Imaging and
599 Optically Manipulating Neuronal Ensembles. *Annu. Rev. Biophys.* 46, 271-293 (2017).
- 600 32. Lecoq,J., Orlova,N., & Grewe,B.F. Wide. Fast. Deep: Recent Advances in
601 Multiphoton Microscopy of In Vivo Neuronal Activity. *J. Neurosci.* 39, 9042-9052 (2019).
- 602 33. Keemink,S.W. *et al.* FISSA: A neuropil decontamination toolbox for calcium imaging
603 signals. *Sci. Rep.* 8, 3493 (2018).
- 604 34. Lin,T.Y., Goyal,P., Girshick,R., He,K., & Dollar,P. Focal Loss for Dense Object
605 Detection. *IEEE Trans. Pattern. Anal. Mach. Intell.* 42, 318-327 (2020).
- 606 35. Dubbs,A., Guevara,J., & Yuste,R. moco: Fast Motion Correction for Calcium
607 Imaging. *Front Neuroinform.* 10, 6 (2016).
- 608 36. Weiss,K., Khoshgoftaar,T.M., & Wang.D.D. A survey of transfer learning. *Journal of*
609 *Big Data* 3, (2016).
- 610 37. Cogswell,M., Ahmed,F., Girshick,R., Zitnick,L., & Batra,D. Reducing overfitting in
611 deep networks by decorrelating representations. *ICLR Conference Track*
612 *Proceedings*(2016).
- 613 38. Charles,A.S., Song,A., Gauthier,J.L., Pillow,J.W., & Tank,D.W. Neuronal Anatomy
614 and Optical Microscopy (NAOMi) Simulation for evaluating calcium imaging methods.
615 *BioRxiv*(2019).
- 616 39. Cicek,O., Abdulkadir,A., Lienkamp,S.S., Brox,T., & Ronneberger,O. 3D U-net:
617 learning dense columetric segmentation from sparse annotation. *Proceedings of the*

618 *International Conference on Medical Image Computing and Computer-Assisted*
619 *Intervention*(2016).

- 620 40. Goard,M. & Dan,Y. Basal forebrain activation enhances cortical coding of natural
621 scenes. *Nat. Neurosci.* 12, 1444-1449 (2009).
- 622 41. Vinck,M., Batista-Brito,R., Knoblich,U., & Cardin,J.A. Arousal and locomotion make
623 distinct contributions to cortical activity patterns and visual encoding. *Neuron* 86, 740-754
624 (2015).
- 625 42. Jacobs,E.A.K., Steinmetz,N.A., Peters,A.J., Carandini,M., & Harris,K.D. Cortical
626 State Fluctuations during Sensory Decision Making. *Curr. Biol.*(2020).
- 627 43. Grawe,B.F. *et al.* Neural ensemble dynamics underlying a long-term associative
628 memory. *Nature* 543, 670-675 (2017).
- 629 44. Shuman,T. *et al.* Breakdown of spatial coding and interneuron synchronization in
630 epileptic mice. *Nat. Neurosci.* 23, 229-238 (2020).
- 631 45. Villette,V. *et al.* Ultrafast Two-Photon Imaging of a High-Gain Voltage Indicator in
632 Awake Behaving Mice. *Cell* 179, 1590-1608 (2019).
- 633 46. Weisenburger,S. *et al.* Volumetric Ca(2+) Imaging in the Mouse Brain Using Hybrid
634 Multiplexed Sculpted Light Microscopy. *Cell* 177, 1050-1066 (2019).
- 635 47. Pluta,S.R., Telian,G.I., Naka,A., & Adesnik,H. Superficial Layers Suppress the Deep
636 Layers to Fine-tune Cortical Coding. *J. Neurosci.* 39, 2052-2064 (2019).
- 637 48. Antonini,A. *et al.* Extended field-of-view ultrathin microendoscopes for high-
638 resolution two-photon imaging with minimal invasiveness. *Elife.* 9, (2020).
- 639 49. Brondi,M. *et al.* High-Accuracy Detection of Neuronal Ensemble Activity in Two-
640 Photon Functional Microscopy Using Smart Line Scanning. *Cell Rep.* 30, 2567-2580
641 (2020).

- 642 50. Forli,A. *et al.* Two-Photon Bidirectional Control and Imaging of Neuronal Excitability
643 with High Spatial Resolution In Vivo. *Cell Rep.* 22, 3087-3098 (2018).
- 644 51. Vecchia,D. *et al.* Temporal Sharpening of Sensory Responses by Layer V in the Mouse
645 Primary Somatosensory Cortex. *Curr. Biol.*(2020).
- 646 52. Mori,T. *et al.* Inducible gene deletion in astroglia and radial glia--a valuable tool for
647 functional and lineage analysis. *Glia* 54, 21-34 (2006).
- 648 53. Castello-Waldow,T.P. *et al.* Hippocampal neurons with stable excitatory connectivity
649 become part of neuronal representations. *PLoS. Biol.* 18, e3000928 (2020).
- 650 54. Dombeck,D.A., Harvey,C.D., Tian,L., Looger,L.L., & Tank,D.W. Functional imaging
651 of hippocampal place cells at cellular resolution during virtual navigation. *Nat. Neurosci.*
652 13, 1433-1440 (2010).
- 653 55. Moretti,C., Antonini,A., Bovetti,S., Liberale,C., & Fellin,T. Scanless functional
654 imaging of hippocampal networks using patterned two-photon illumination through GRIN
655 lenses. *Biomed. Opt. Express* 7, 3958-3967 (2016).
- 656 56. He,K., Zhang,X., Ren,S., & Sun,J. Deep residual learning for image recognition.
657 *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern
658 Recognition*(2016).
- 659 57. Kingma,D.P. & Ba,J.L. Adam: a method for stochastic optimization. *arXiv*(2015).
- 660 58. Herlihy,M. & Shavit,N. *The art of multiprocessor programming*(Elsevier,2006).
- 661 59. Malisiewicz,T., Gupta,A., & Efros,A.A. Ensemble of exemplar-SVMs for object
662 detection and beyond. *Proceedings of the IEEE International Conference on Computer
663 Vision*(2011).
- 664 60. Abadi,M. *et al.* TensorFlow: a system for large-scale machine learning. *arXiv*(2016).
- 665 61. Van Rossum,G. & Drake,F.L. *Python 3 Reference Manual*(CreateSpace,2009).

666 62. Miller,J.E., Ayzenshtat,I., Carrillo-Reid,L., & Yuste,R. Visual stimuli recruit
667 intrinsically generated cortical ensembles. *Proc. Natl. Acad. Sci. U. S. A* 111, E4053-E4061
668 (2014).

669 **Methods**

670 *Animals*

671 All experiments were carried out in accordance with the guidelines of the European Communities
672 Council Directive, and were approved by the National Council on Animal Care of the Italian Ministry
673 of Health (authorization #34/2015-PR, #1134/2015-PR, and #61/2019-PR).

674

675 Wild type (wt) C57BL/6J mice were purchased from Charles River Laboratories (Calco, Italy; strain
676 code #632), transgenic *Scnn1a-cre* (B6;C3-Tg(Scnn1a-cre)3Aibs/J; JAX #009613), and *Ai95D*
677 (B6;129S-Gt(ROSA)26Sor^{tm1(CAG-GCaMP6f)Hze}/J; JAX #024105) were purchased from Jackson
678 Laboratories (Bar Harbor, USA). *Scnn1a-cre* mice express the enzyme Cre in a subpopulation of layer
679 IV neurons⁴⁷ and of VPM cells⁴⁸. Animals were housed in individually ventilated cages under a 12-
680 hr light:dark cycle. A maximum of five animals per cage was allowed with *ad libitum* access to food
681 and water. Mice of both sexes were used for experiments.

682

683 *Viral injections and surgical procedures*

684 We expressed GCaMP6 or GCaMP7 through the following viral vectors
685 AAV1.Syn.Flex.GCaMP6s.WPRE.SV40 (Addgene viral prep # 100845-AAV1),
686 AAV1.Syn.Flex.GCaMP6f.WPRE.SV40 (Addgene viral prep # 100833-AAV1) purchased from the
687 University of Pennsylvania Viral Vector Core, or AAV1.Syn.Flex.GCaMP7f.WPRE.SV40 (Addgene
688 viral prep #104492-AAV1) purchased from Addgene. For CA1 imaging, we expressed jRCaMP1a
689 using co-injection of AAV1.CAG.Flex.NES-jRCaMP1a.WPRE.SV40 (Addgene viral prep # 100846-
690 AAV1) and AAV1.CamKII 0.4.Cre.SV40 (Addgene viral prep # 105558-AAV1) purchased from the
691 University of Pennsylvania Viral Vector Core.

692

693 For LIV imaging, we used a total of 29 mice. Specifically, 17 *Scnn1a-cre* mice injected with a viral
694 vector transducing GCaMP6s, 6 *Scnn1a-cre* mice injected with a virus transducing GCaMP6f, 3
695 *Scnn1a-cre* mice injected with a virus transducing GCaMP7s, and 3 *Scnn1a-cre* crossed with *Ai95D*
696 mice. Mice between post-natal days 30 and 33 were anesthetized with 2% isoflurane (IsoFlu, Zoetis,
697 IT) in 0.8 % oxygen, placed into a stereotaxic apparatus (Stoelting Co, Wood Dale, IL), and
698 maintained on a warm platform at 37°C for the whole duration of the anesthesia. Viral injection in
699 mice used for LIV imaging was carried out similarly to ⁴⁹ and ⁵⁰. Briefly, a scalp incision was
700 performed after local administration of lidocaine (2 %) and then two small holes were drilled on the
701 skull above the right/left somatosensory cortex at 1.2 mm and 2 mm posterior (P) to the bregma
702 suture, 2.8 mm and 3 mm lateral (L) to the sagittal sinus. A micropipette was slowly inserted into the
703 cortical tissue until the tip reached a depth of 0.3 mm below the pia mater. 200 nL of GCaMP6 virus
704 were injected at 20 - 60 nl/min by means of a hydraulic injection apparatus driven by a syringe pump
705 (UltraMicroPump, WPI, Sarasota, FL). The pipette was then further lowered to reach 0.4 mm below
706 the *pia mater*, and a second injection was performed. This procedure was repeated for the second
707 injection site. The injected solution contained 10¹² viral genomes/ml diluted 1:1 in artificial cerebro-
708 spinal fluid (aCSF: 127 mM NaCl, 3.2 mM KCl, 2 mM CaCl₂, 1 mM MgCl₂ and 10 mM HEPES, pH
709 7.4). The exposed skull was then cleaned, and the skin sutured and cleansed with Iodopovidone
710 (Betadine®, Meda pharma, Milan, Italy). Mice were monitored until full recovery from the anesthesia.
711 In mice used for imaging in awake conditions, a custom metal bar was sealed to the skull using
712 Vetbond (3 M, St. Paul, MN, USA) and dental cement (Paladur, Kulzer GmbH, Hanau, Germany).
713 The exposed bone was covered using the silicone elastomer KWIK-Cast (World Precision
714 Instruments, Friedberg, DE) and an intraperitoneal injection of antibiotic (BAYTRIL, Bayer, DE) was
715 performed. 2-4 weeks after virus injection, mice used for imaging in LIV during anesthesia were
716 injected with urethane (1.65 g/kg, 16.5 % in saline solution, i-p.). A scalp incision was performed

717 after local administration of Lidocaine (2 %). A circular craniotomy was opened over the
718 somatosensory cortex, in the area where green fluorescence was clearly visible using an
719 epifluorescent microscope. The surface of the brain was kept moist with aCSF. A heating pad
720 underneath the animal was set at 35.5-37°C. Respiration rate, eyelid reflex, vibrissae movements, and
721 reactions to tail pinching were monitored throughout the surgery. Mice were then moved under the
722 two-photon microscope, kept at 37 °C with a heating pad, and the brain surface irrigated with aCSF.
723 Imaging began one hour after the end of the surgery. Before imaging LIV activity in awake animals,
724 mice were habituated to head-fixation similarly to ⁵¹. In brief, habituation lasted for 7 to 10 days,
725 during which they were head restrained for increasing periods (from 15 minutes to one hour), while
726 running or standing on a custom made treadmill. On the day of the experiment, the habituated mouse
727 was anesthetized with a mixture of isoflurane and oxygen (0.8 – 2 %), and a craniotomy performed
728 similarly to what described above. After surgery, the animal was head fixed and allowed to recover
729 under the microscope for at least one hour before imaging.

730

731 For VPM imaging, we used a total of 4 mice. Viral injections and aberration-corrected
732 microendoscopes insertion in mice used for VPM imaging were performed in Scnn1a-Cre mice as in
733 ⁴⁸. Mice were anesthetized as previously described. A single craniotomy was performed at stereotaxic
734 coordinates P 1.7 mm, L 1.6 mm. A micropipette was lowered to a depth of 3 mm below the brain
735 surface. 0.5 - 1 µl of GCaMP6s virus containing solution (containing 10¹² viral genomes/ml diluted
736 1:4 in aCSF) were injected at 30 - 50 nl/min. Following virus injection, a craniotomy (area: 600 µm
737 x 600 µm) was performed at stereotaxic coordinates P 2.3 mm, L 2 mm. A thin column of brain tissue
738 was displaced with a glass cannula (ID = 300 µm, OD = 500 µm; Vitrotubs, Vitrocom Inc., Mounting
739 Lakes, NJ) and a microendoscope was slowly inserted into the cannula track using a custom holder,

740 down to 3 mm from the brain surface. The microendoscope was finally secured by acrylic adhesive
741 and dental cement. Imaging was performed 2-4 weeks after endoscope implantation.

742

743 For CA1 imaging, we used a total of 2 mice. Before surgery, mice were medicated with an
744 intramuscular bolus of Dexamethasone (Dexadreson, 4 gr/kg). After scalp incision, a 0.5 mm
745 craniotomy was drilled at stereotaxic coordinates P 1.75 mm, L 1.35 mm. A micropipette was lowered
746 1.40 mm below the brain surface. 0.8 μ l of viral solution (containing a mixture of CamKII-Cre and
747 jRCaMP1a viruses at 10^{12} viral genomes/ml diluted 1:1 in aCSF) was injected at 100 nL/min in *Ai95*
748 crossed with *Glast-cre-ERT2* mice ⁵². Inducible Glast-cre was not activated after viral injection,
749 resulting in CA1 neuronal staining with jRCaMP1a and GCaMP6f. A stainless-steel screw was
750 attached to the skull, and a chronic hippocampal window was implanted as described in ^{53,54}. A 3 mm
751 craniotomy was opened, centered at coordinates P 2.00 mm, L1.80 mm. The *dura mater* was removed
752 using fine forceps, and the cortical tissue overlaying the hippocampus slowly aspirated using a blunt
753 needle coupled to a vacuum pump. During aspiration, the exposed tissue was continuously irrigated
754 with aCSF. Aspiration was stopped once the fibers of the external capsule were exposed. A cylindrical
755 optical window made of a thin-walled stainless-steel cannula (OD, 3 mm; ID, 2.77 mm; height, 1.50
756 - 1.60 mm) attached to a 3 mm diameter round coverslip, was fitted to the craniotomy in contact to
757 the external capsule. A thin layer of silicone elastomer was used to surround the interface between
758 the brain tissue and the steel surface of the optical window. A custom stainless-steel headplate was
759 attached to the skull using epoxy glue. All the components were finally fixed in place using black
760 dental cement and the scalp incision was sutured to adhere to the implant. All the animals received
761 an intraperitoneal bolus of antibiotic (BAYTRIL, Bayer, DE) to prevent postsurgical infections.

762

763 *Functional two-photon imaging*

764 Two-photon imaging was performed using a chameleon ultra II pulsed laser (80 MHz pulse frequency,
765 Coherent Inc, Santa Clara, CA, USA) tuned at 920 nm for GCaMP6/7 imaging and at 990 nm for
766 dual color imaging. Excitation power was 30 - 110 mW as measured under the microscope objective
767 and controlled via a Pockel cell (Conoptics Inc, Danbury CT, USA.). An Ultima II scanhead (Bruker
768 Corporation, Milan, Italy) equipped with 3 mm raster scanning galvanometers (6215H, Cambridge
769 Technology, Bedford, MA) and standard photomultiplier tubes (Hamamatsu Photonics, Tokyo, Japan)
770 and an Ultima Investigator (Bruker Corporation, Milan, Italy), equipped with 6 mm raster scanning
771 galvanometers, movable objective mount, and multi-alkali photomultiplier tubes were used. The three
772 objectives were: 25x/1.05 N.A. (Olympus Corp., Tokyo, JP) for LIV imaging, 20x/0.5NA (Zeiss,
773 Oberkochen, Germany) for VPM endoscopic imaging, and 16x/0.8 N.A. (Nikon, Tokyo, Japan) for
774 CA1 experiments.

775

776 For LIV imaging, dwell time was 4.4 μ s, photomultiplier voltage was 777 V, zoom factor was always
777 1, pixel size was 0.77 μ m, acquisition frame rate ranged between 0.5 – 3 Hz for a 512 pixels x 512
778 pixels image. Fluorescence values spanned 95 % of the available dynamic range (16 bit). For dual
779 color CA1 imaging, pixel dwell-time was set at 4 μ s, photomultiplier voltage was 777 V, zoom factor
780 was always 5, pixel size was 0.634 μ m, acquisition frame rate was 3.03 Hz for a 256 pixels x 256
781 pixels image. For VPM imaging, the set up was similar to the one described in ^{48, 55}, pixel dwell-time
782 was set at 4 μ s, photomultiplier voltage was 810 V, zoom factor was always 1, pixel size was 2.19
783 μ m, acquisition frame rate was 2.66 Hz for a 196 pixels x 196 pixels image. Imaging sessions lasted
784 1 hour for CA1, VPM, and awake LIV experiments. They lasted 4 hours for the anesthetized LIV
785 condition. After awake imaging sessions, animals were returned to their home cages.

786

787 *Training and validation datasets*

788 In the absence of a generally accepted wide-scale annotated dataset of two-photon calcium imaging,
789 we built a dataset of *in vivo* t-series collected using raster scanning acquisitions. A total of 197 t-series
790 (average frame number per time series: 597 ± 262 , average frame rate: 2.3 ± 1.5 Hz) were included
791 in the dataset: 76 t-series from CA1 imaging of principal neurons stained with both GCaMP6f (38 t-
792 series) and jRCaMP1a (38 t-series); 121 t-series from cortical LIV imaging of principal neurons
793 stained with virally injected GCaMP6s (107 t-series), GCaMP6f (4 t-series), GCaMP7f (5 t-series),
794 GCaMP6f expressed in transgenic animals (Scnn1a-cre x Ai95; 5 t-series). Training and validation
795 datasets contained 160 (118 from LIV, 21 from CA1 GCaMP6f, and 21 from CA1 jRCaMP1a) and
796 37 t-series (13 from LIV, 12 from CA1 GCaMP6f, and 12 from CA1 jRCaMP1a), respectively. To
797 avoid data leakage between training and validation datasets, we manually split t-series including
798 different FOVs in the datasets.

799

800 *Additional datasets*

801 Four additional datasets were selected and used for validation purposes only:
802 1) VPM microendoscopic imaging t-series in awake head restrained mice (9 t-series).
803 2) The publicly available Allen Brain Observatory (ABO) visual coding dataset (19 t-series,
804 <https://observatory.brain-map.org/visualcoding>). T-series identification numbers: 501271265,
805 501484643, 501574836, 501704220, 501729039, 501836392, 502115959, 502205092, 502608215,
806 503109347, 504637623, 510214538, 510514474, 510517131, 527048992, 531006860, 539670003,
807 540684467, 545446482.
808 3) The publicly available Neurofinder (NF) challenge dataset (28 t-series,
809 <https://github.com/codeneuro/neurofinder>). T-series identification numbers: neurofinder.00.00,

810 neurofinder.00.01, neurofinder.00.02, neurofinder.00.03, neurofinder.00.04, neurofinder.00.05,
811 neurofinder.00.06, neurofinder.00.07, neurofinder.00.08, neurofinder.00.09, neurofinder.00.10,
812 neurofinder.00.11, neurofinder.01.00, neurofinder.01.01, neurofinder.02.00, neurofinder.02.01,
813 neurofinder.03.00, neurofinder.04.00, neurofinder.04.01, neurofinder.00.00.test,
814 neurofinder.00.01.test, neurofinder.01.00.test, neurofinder.01.01.test, neurofinder.02.00.test,
815 neurofinder.02.01.test, neurofinder.03.00.test, neurofinder.04.00.test, neurofinder.04.01.test.

816 4) A single t-series of mesoscopic full field imaging from ⁷.

817 No preprocessing was performed on the VPM, ABO, NF, and mesoscopic t-series. All t-series were
818 manually annotated *de novo* by the two graders working independently. The consensus ground truth
819 was obtained as described for the training and validation dataset (see also Supplementary Table 1-2).

820

821 *Image processing and consensus labelling*

822 The CITE-On image detector was based on purely morphological features extracted from imaging
823 data. No information from the dynamic fluorescence signal in the t-series was used to detect putative
824 cells. Each imaging t-series was corrected for lateral displacements using MOCO ³⁵. T-series were
825 aligned using the first frame as reference, without downsampling and with a maximum possible shift
826 of 13 pixels. The 8-bit median projection of each t-series was then computed on the motion-corrected
827 t-series. The resulting images (one *per* t-series) were globally sharpened to better visualize cell
828 shapes. A gamma correction of 0.3 was applied, and the dynamic range was linearly adjusted
829 normalizing across the whole 8 bit depth. Processed images were named “enhanced median
830 projections” (EMPs) and were used to define our GT labelling. Two graders independently labelled
831 each EMP. LabelImg (<http://github.com/tzutalin/labelImg>) was used to define a single object class by
832 manually drawing bounding boxes around every visible cell soma in the EMP. The surface of each
833 bounding box was manually defined in order to tightly surround the cell shape. Boxes were allowed

834 to overlap. Coordinates and surface of each bounding box for all EMPs were saved in a standard VOC
835 format where each file reported the top left and bottom right coordinates (in pixels) for each bounding
836 box. The GT was defined as the union between the set of notations from the two graders.

837

838 *Image detector training*

839 CITE-On is based on a fully convoluted single-shot image detector, RetinaNet³⁴. Briefly, a feature
840 pyramid network was constructed from residual layers of the ResNet50 feature extractor⁵⁶. This
841 feature pyramid was then fed to two separate sets of convolution filters: one computing the label score
842 (classification subnet), the second performing bounding box regression from priors (regression
843 subnet). We used the open-source Keras implementation of this architecture
844 (<http://github.com/fizyr/keras-retinanet>). Given the relatively small set of training data available, we
845 could not perform *ab initio* training. Rather, we employed a transfer learning approach. Starting from
846 the network trained on a large-scale dataset of natural images, we fine-tuned the weights of the
847 classification and regression subnets, while freezing the weights of the feature pyramids. We used
848 “plain” median projections obtained from the motion corrected t-series and linearly normalized across
849 the bit range. The resulting projections were then upsampled in order to obtain images were the short
850 side was 800 pixels long, while the long side did not exceed 1,333 pixels. Since the input layer of the
851 network accepted a three-channel image, the same image was replicated for each channel without
852 changing any parameter. These last two image conversions were necessary as the network was
853 originally trained on RGB images, and the transfer learning approach did not allow input
854 modifications to the image shape. The network was trained with a regression L1 loss function (Mean
855 Absolute Errors, MAE, <https://rishi.github.io/ml/2015/07/28/l1-vs-l2-loss/>), with focal loss
856 (<http://arxiv.org/abs/1708.02002>) and the Adam optimizer⁵⁷ with learning rate 10-5 and clipnorm 10-
857 3 (<http://github.com/keras-team/keras/issues/510>)³⁴ modified by reducing the learning rate on loss

858 plateau with a factor of 0.1. The network was trained for 17 epochs, each consisting of 1000 training
859 steps of batch size 1.

860

861 *CITE-On offline pipeline*
862 Two-photon calcium imaging t-series were first corrected for lateral artifacts using MOCO (as
863 described above). The median projections were then computed, normalized, upscaled to the target
864 input size, and converted to 8 bit RGB images. The resulting images were fed to the image detection
865 network. Upscaling factor and score threshold were the only two parameters defined *a priori*.
866 Upscaling factor was adjusted in order to have the smallest feature in each image inscribed in a 32
867 pixels² box. This was because the smallest prior box encoded in the network was 32 pixels². To avoid
868 using a shallower convolutional residue that would have carried less information, we decided to adjust
869 the up-scaling factor instead of reducing the size of the smallest prior box. After exploring the
870 parameter space on different datasets, we defined an up-scaling factor of 1 for the training and
871 validation datasets. The optimal up-scaling factors were 2.3 for the ABO dataset, and between 1.7
872 and 3.1 for the NF datasets (these datasets were heterogeneous in terms of FOV size and pixel size).
873 In the VPM dataset, image projections were characterized by an altered magnification factor in the
874 radial direction, due to the optical properties of the corrected microendoscopes used for imaging⁴⁸.
875 These images went through a further preprocessing step to correct for this field distortion. Up-scaling
876 factor for the VPM dataset was 1.4. Each bounding box was associated with a score, representing
877 network confidence in cell detection. Bounding boxes with intersection over union (IoU) < 20 % were
878 considered as separate neuronal identities. When IoU of two bounding boxes was > 20 %, the
879 bounding box with the highest score was retained. Results of the image detector were filtered by
880 applying a threshold on the output score provided by the network and optimized for each dataset.

881

882 *CITE-On online pipeline*

883 In the online pipeline, individual raw imaging frames were continuously grabbed from a streaming
884 source (e.g. live microscope output) and processed. To simulate this process, we individually
885 imported in the CITE-On pipeline each frame of each raw t-series. Single frames were passed on to
886 the trace extractor and to a buffer. The buffer stored the number of frames sufficient to produce an
887 average projection. Once the buffer was filled, the projection was computed, sent to the image
888 detector, and the buffer emptied.

889

890 Detections were performed using the same procedure described for the offline pipeline. Detection
891 results were fed to a custom tracking algorithm detecting all the overlaps between current and
892 previous detections, and designed in order to maximize overlap between putative matching boxes.
893 For every detection matching a previous one, the coordinates of the relative bounding box were
894 updated to the last one. For each new detection having an IoU < 0.25 % with all the previous
895 detections, a new identity was created. In case of identities not actively detected in the current frame,
896 relative coordinates were anyway updated using a rigid shift calculated as the mean shift obtained
897 from the active identities. In this way, we aimed to minimize the effect of motion artifacts and identity
898 switch without implementing online motion correction approaches. A simple dynamic segmentation
899 was then performed for each identity. At each raw frame, the interval between the 80th and the 95th
900 percentile of the pixel fluorescence intensity distribution inside each bounding box was averaged to
901 extract the raw functional trace. At each frame, background signal corresponded to the average
902 fluorescence of all the pixels in the FOV not belonging to any bounding box. This frame-wise value
903 was subtracted from all the individual raw functional traces. In order to optimize real-time
904 performance for high frame rate acquisitions (above 3 Hz, including all ABO and some NF t-series),
905 the entire pipeline was implemented in a multiprocessing scheme where one process was responsible

906 for data loading, one for image preprocessing, and for sending its output to the CNN detector
907 accelerated over GPU and tracking parts, while the remaining CPU cores were dedicated to real-time
908 trace extraction given the parallel nature of the problem⁵⁸. This implementation scheme allowed for
909 cell detection update (up to 10 Hz) and functional trace extraction update from all identities (100 Hz)
910 to operate as parallel and asynchronous processes.

911

912 *Parameter exploration for the object detector*

913 To find the best operating parameters for the object detector, we quantified offline CITE-On
914 performance while systematically exploring various plausible values of score threshold and upscaling
915 factor. For all datasets, we evaluated the F-1 across a set of net size multipliers between 0.6 and 3.4
916 in steps of 0.2. We also explored score thresholds between 0.05 and 0.95, in steps of 0.05. This
917 mapping strategy allowed us to define the optimal combination of score threshold and upscaling
918 factor for each input data. The upscaling factor was dependent on the ratio between the average box
919 surface and the whole FOV surface, while the score threshold presented non obvious dependence on
920 trivial image statistics. Therefore, we determined the upscaling factor according to the acquisition
921 parameters and the relative score thresholds, in order to maximize the F-1 score for each dataset. For
922 the online pipeline, we used appropriate upscaling factors and proceeded by exploring the dependency
923 of F-1 from the score threshold and from the number of averaged frames in each detection.

924

925 *Trace extraction: comparison between CITE-On and other methods*

926 We compared CITE-On trace extraction with trace extraction performed with a popular state-of-the-
927 art method based on CNMF, CaImAn¹⁰. Briefly, we provided binary masks corresponding to the
928 CITE-On detected bounding boxes and used these masks to seed the CNMF algorithm. Seeded-
929 CNMF first calculated the temporal background component using pixels that were not included in

930 any mask. We compared this background component to the CITE-On background traces used for trace
931 correction. The subsequent step of the seeded-CNMF algorithm estimated temporal components and
932 spatial footprints, constrained to be non-zero only at the location of the binary masks. Using this
933 strategy, we obtained fluorescent traces from putative cells detected in the same locations as those
934 detected by CITE-On, allowing for a one-to-one trace comparison between algorithms. It is important
935 to note that, for this analysis, we set the order of the autoregressive model of the CNMF to zero
936 because we were not interested in trace deconvolution, but only in correcting for background
937 contamination.

938

939 *Local vs. global background signal correlation*

940 To compare local and global background noise contributions, we used the same approach for
941 background noise subtraction, but considering only the pixels in the vicinity of each cell. The vicinity
942 of a cell was defined as all the pixels in a concentric rectangular box double the size of the box
943 detected by CITE-On, with no overlap with other bounding boxes. We then calculated the cross
944 correlation at lag zero between the local noise for each cell and the global background trace.

945

946 *Tiled detection on mesoscopic images*

947 For large scale datasets ($> 3000 \text{ pixel}^2$) as the mesoscopic imaging dataset⁷, requiring large amounts
948 of GPU memory ($> 500 \text{ GB}$), we implemented a tiled detection approach. We divided each
949 mesoscopic FOV in a configurable number of tiles with configurable overlap factor, in order to batch
950 process all tiles up to the limit of the available GPU memory. Once all detections were computed,
951 they were appropriately shifted back to the original position in the FOV and a Non Maximum
952 Suppression⁵⁹ was performed in order to remove duplicate boxes in regions of the FOV where overlap
953 between tiles occurred.

954

955 *Performance and metrics*

956 Performance of the CITE-On object detector was calculated employing standard Precision, Recall
957 and F-1 metrics. In order to determine true or false positives, we computed the IoU between each
958 ground truth and detected boxes. A custom algorithm was used to match the two sets using a cut-off
959 threshold of 0.5.

960

961 *Hardware and software for data analysis*

962 All the data analysis procedures presented in this work were performed on a Dell Precision 7920
963 desktop with an Intel Xeon Silver 4110 @ 2.1 GHz 8 core CPU, 32 GB DDR4-2666 ECC RAM,
964 NVIDIA Quadro RTX5000 GPU, 512 GB NVMe SSD, and 2 TB 7200 rpm HDD.

965

966 All processing steps, including network training and validation, were carried out under
967 Keras/Tensorflow software libraries⁶⁰. Image processing and data analysis were carried out using
968 Python Language Reference 3⁶¹.

969

970 *Statistics*

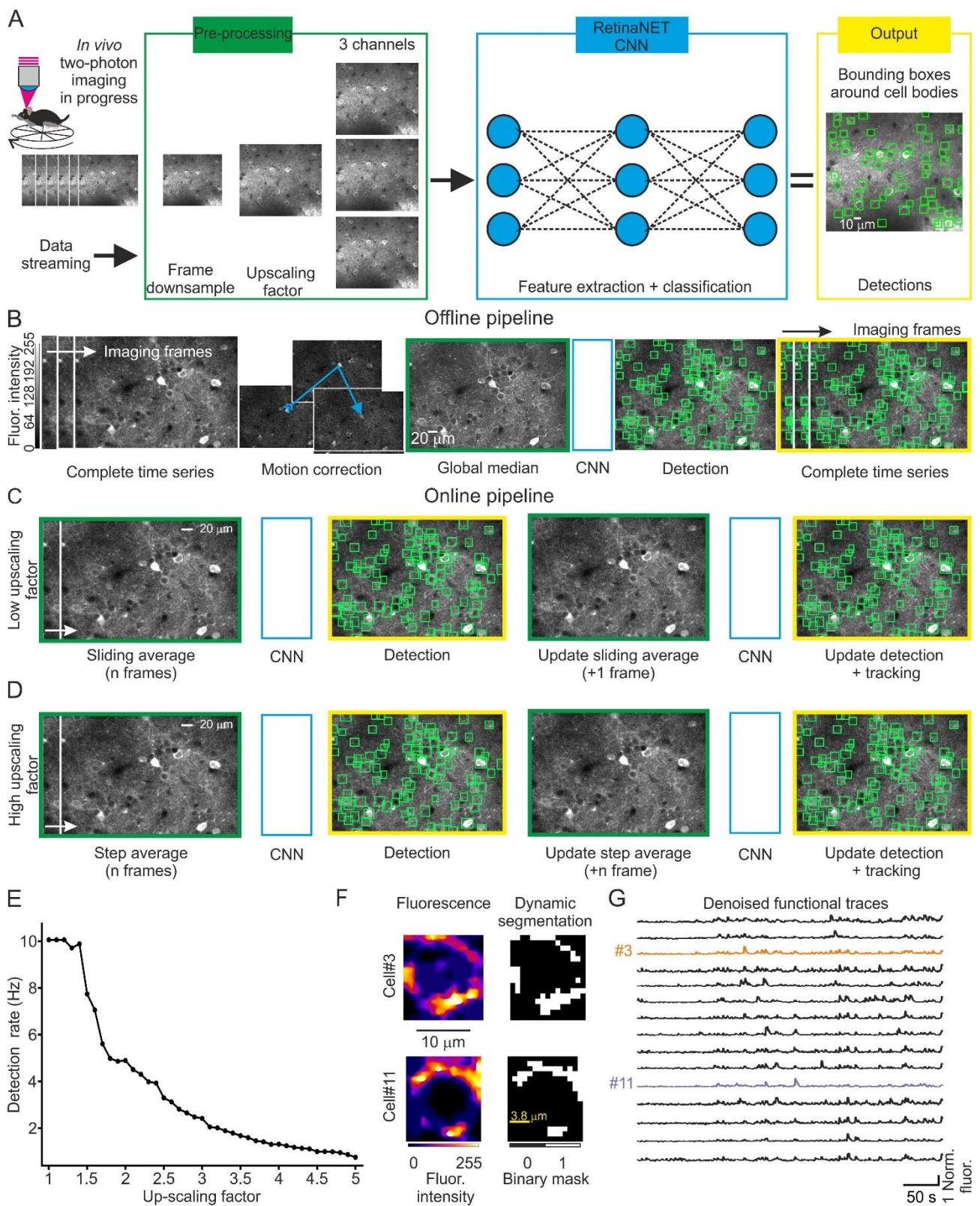
971 Values were expressed as mean \pm sd, unless otherwise stated. The number of samples (N) and p values
972 are reported in the figure legends or in the text. No statistical methods were used to pre-determine
973 sample size, but sample size was chosen based on previous studies^{3, 62}. All recordings with no
974 technical issues were included in the analysis and blinding was not used in this study. Statistical
975 analysis was performed with the scientific Python ecosystem (SciPy, NumPy 1.19) under Python 3,
976 Python Software Foundation, Python Language Reference 3 (available at <https://www.python.org>). A
977 Kolmogorov-Smirnov test was run on each experimental sample to test for normality. The

978 significance threshold was always set at 0.05. When comparing two paired populations of non-
979 normally distributed data, a two-sample Kolmogorov-Smirnov test. All tests were two-sided, unless
980 otherwise stated.

981

982 *Code availability*

983 The code is available from the corresponding authors on request and it will be shared upon
984 publication.

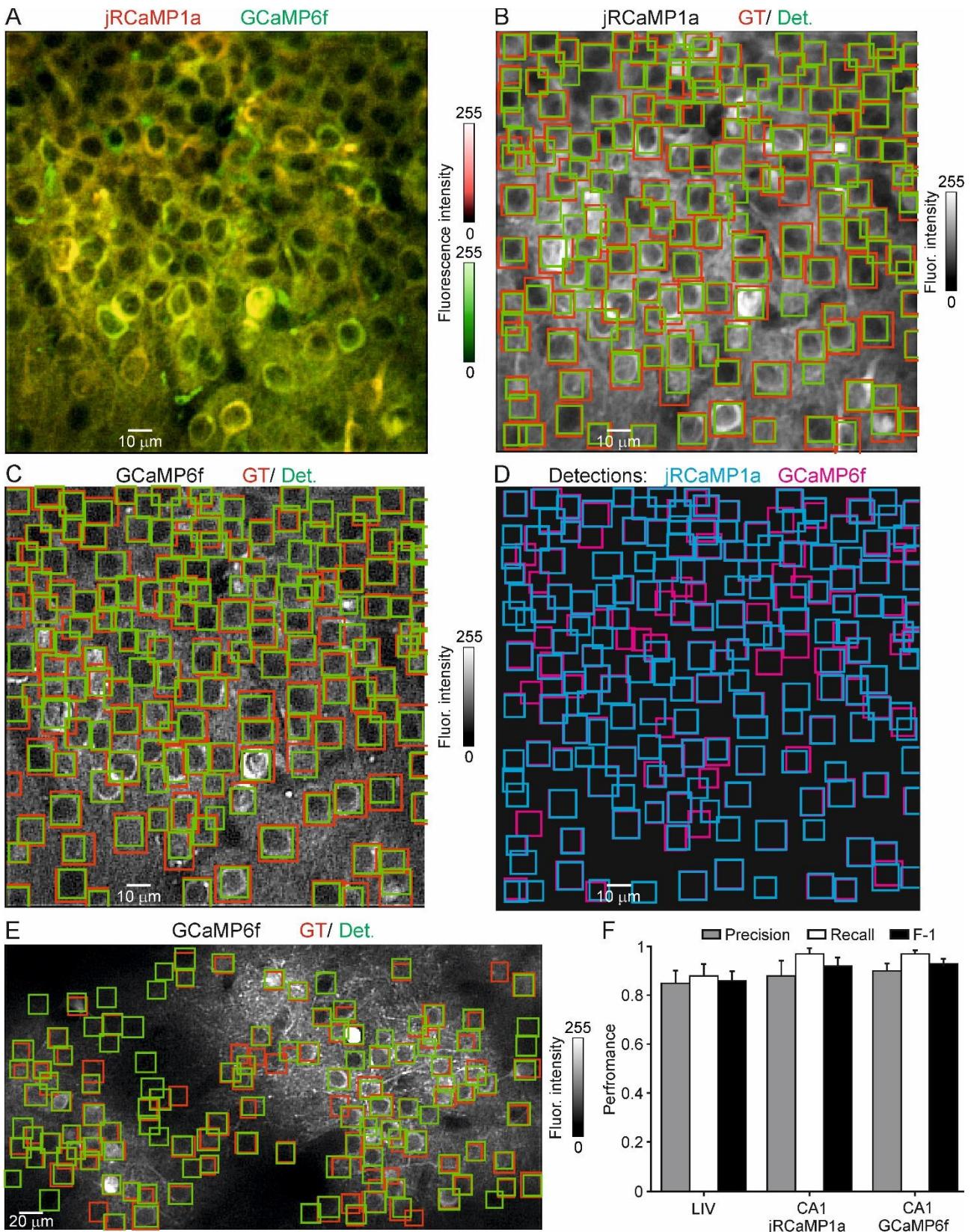


985

Figure 1. Structure and analysis pipeline of CITE-On. A) Schematic of the image detection process in CITE-On. During ongoing two-photon imaging acquisition, individual frames are

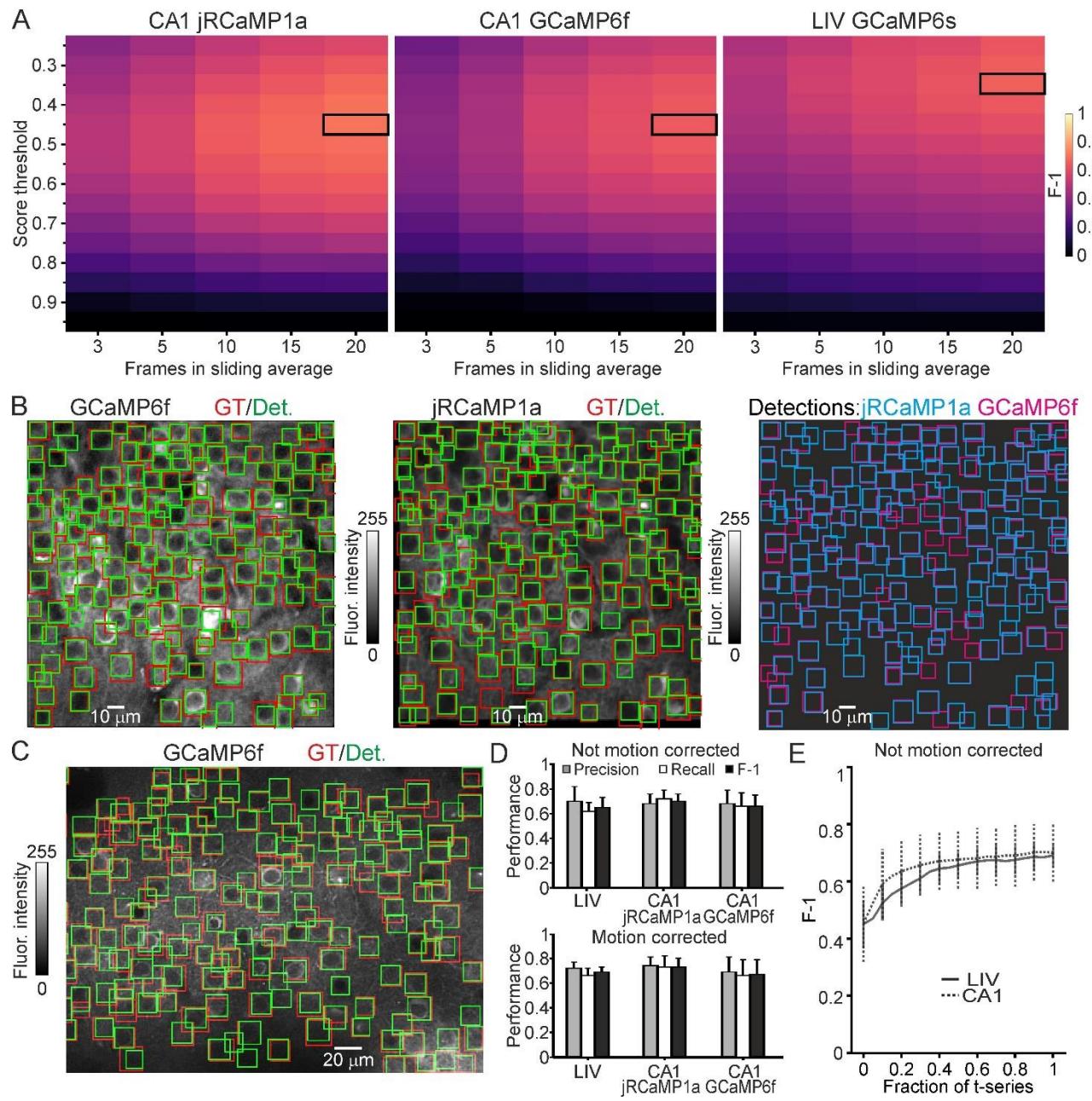
988 transferred to CITE-On as they are completed (left). A pre-processing step (green rectangle) is
989 required ahead of image detection, comprising frame downsampling, image upscaling and triplication
990 of the upscaled image. The result of the preprocessing is then used as input for the CNN (blue
991 rectangle). The CNN output is the detection of neuronal somata in the form of bounding boxes (green
992 squares black and white image on the right). **B)** CITE-On offline pipeline starts with the complete t-
993 series and the correction of motion artifacts (blue arrow, motion correction). Frame downsampling is
994 performed by computing the global median projection of the t-series. The upscaled and triplicated
995 global median (green) is fed to the CNN (blue), a single detection is performed, and the bounding
996 boxes (detection, green squares) are projected onto each frame of the complete t-series (yellow). **C)**
997 In the online pipeline, for data requiring small upscaling factors, a sliding average projection of the
998 first n frames of the ongoing t-series is calculated in the frame downsampling pre-processing step
999 (green). This image is upscaled and triplicated, processed by the CNN (blue), producing the first
1000 detection (yellow). As the next frame of the t-series is acquired, a new sliding average is computed,
1001 again on n frames, but starting from the second frame of the acquisition and including the $n+1^{th}$ one.
1002 The CNN processes this image, updating the detections and starting the tracking system (yellow).
1003 The grayscale shown in this panel applies to all grayscale images in this figure. **D)** For data requiring
1004 high upscaling factors, the pipeline is similar to that in (C), but instead of a sliding average, a step
1005 average is calculated on n frames as the frame downsampling pre-processing step (green). Detections
1006 are updated every n new frames. **E)** Detection rates as a function of the magnitude of the upscaling
1007 factor. Maximum detection rate is 10 Hz for upscaling factor between 1 and 1.5. **F)** Representative
1008 average fluorescence of pixels inside the bounding box relative to two cells (cell #3 and cell #11),
1009 calculated in a single frame of the LIV dataset (GCaMP6s, pseudocolor, left). Associated dynamic
1010 segmentation mask in the same frame (binary mask, right). **G)** Functional traces from $N = 15$

1011 representative cells extracted with online CITE-On pipeline. Traces of cells displayed in (F) are
1012 shown in green.



1013

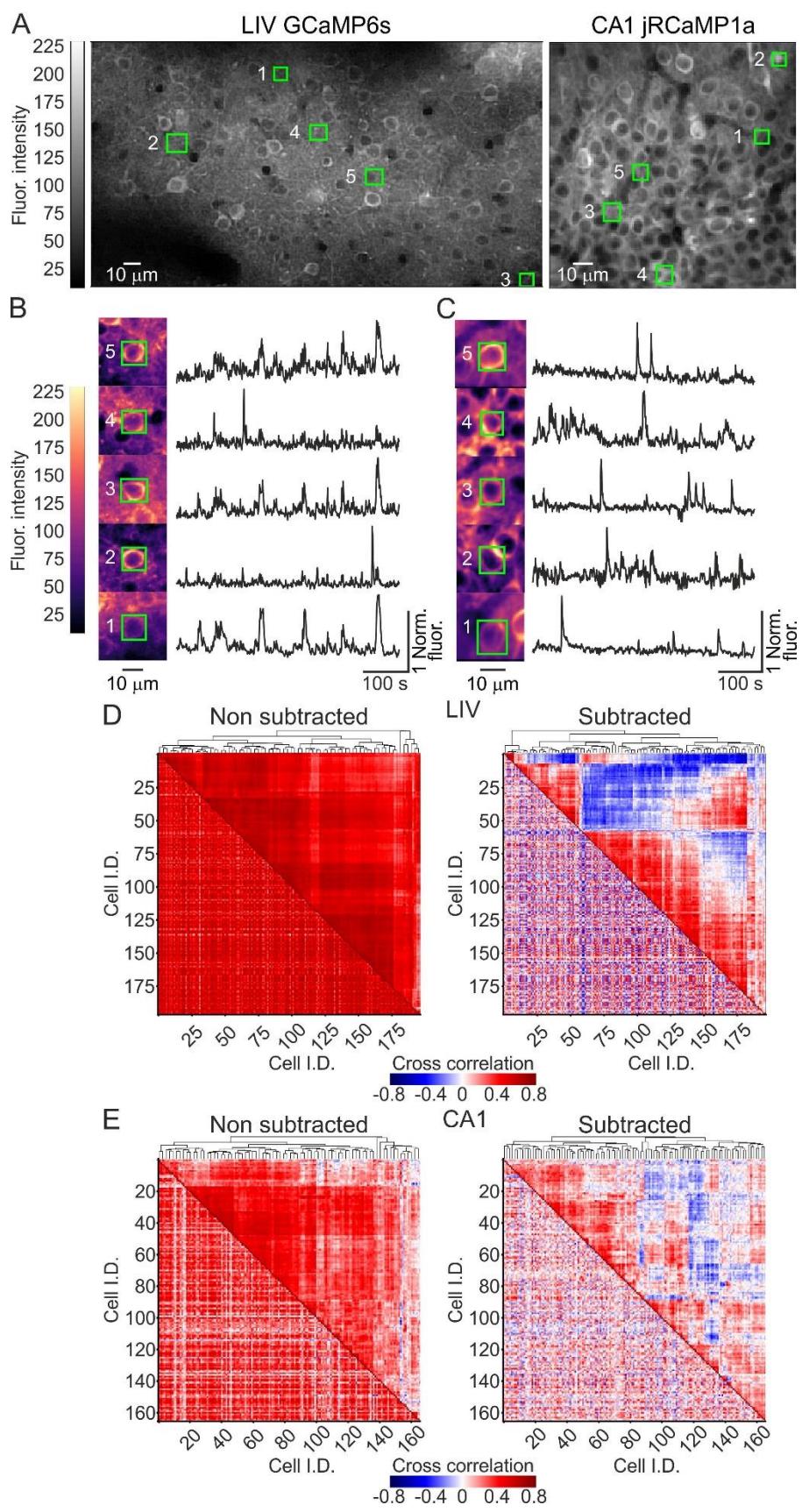
1014 **Figure 2. CITE-On offline performance.** **A)** Representative median projection showing jRCaMP1a
1015 (red) and GCaMP6f (green) expressing CA1 neurons. **B)** GT (red) and CITE-On detections (green)
1016 for the jRCaMP1a channel of the image shown in A. **C)** same as in (B), but for the GCaMP6f channel.
1017 **D)** Superposition of CITE-On detections on jRCaMP1a (cyan) and GCaMP6f (magenta) channels.
1018 **E)** Representative median projection from the LIV dataset with GT (red) and CITE-On detections
1019 (green). **F)** Performance for precision (grey), recall (white), and F-1 (black) obtained with the offline
1020 CITE-On pipeline on the validation t-series of the LIV ($N = 13$), CA1 jRCaMP1a ($N = 12$), and CA1
1021 GCaMP6f ($N = 12$) datasets. In this and in other figures values are expressed as means \pm sd.



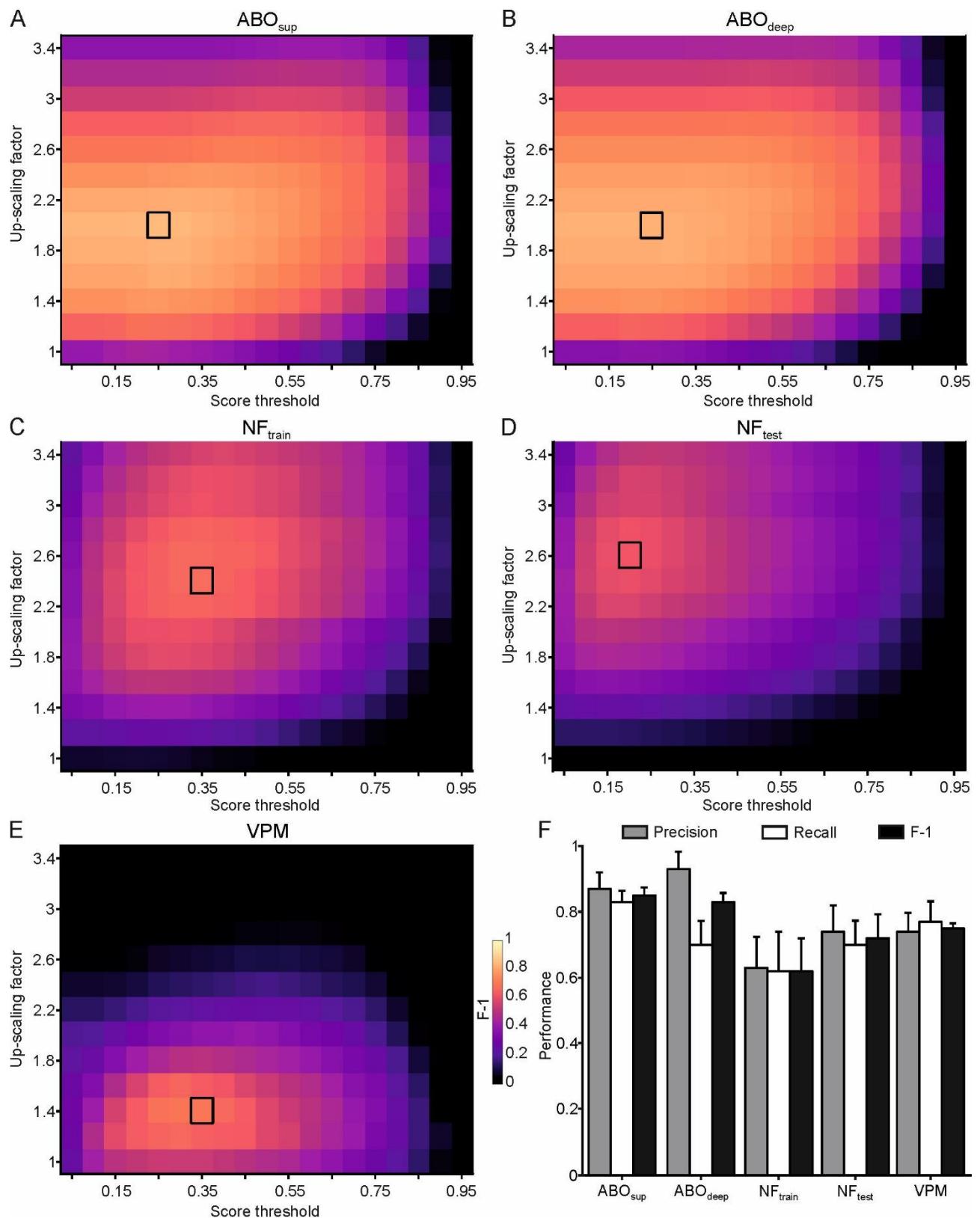
1022

1023 **Figure 3. CITE-On online performance.** **A)** Best parameter search for frame downsampling: F-1
 1024 score (pseudocolor) as a function of score threshold (vertical axis) and number of frames in the sliding
 1025 average (horizontal axis) for LIV (left), CA1 GCaMP6f (middle) and CA1 jRCaMP1a (right). The
 1026 maximal F-1 is indicated with the black rectangle. **B)** Median projections of one representative FOV
 1027 for CA1 jRCaMP1a (left) and one FOV for CA1 GCaMP6f (middle). GT (red) and online detections
 1028 (DET, green) are also shown. In the rightmost panel, the online detections of jRCaMP1a (cyan) and
 1029 GCaMP6f (magenta) are shown. **C)** Same as in (B) but for a representative LIV t-series. **D)** Top:

1030 online performance of Precision (grey), Recall (white), and F-1 (black) for all t-series in the validation
1031 LIV ($N = 13$), CA1 jRCaMP1a ($N = 12$), and CA1 GCaMP6f ($N = 12$) datasets. No motion correction
1032 was performed. Bottom: same as top, but for the motion corrected t-series. Results of Kolmogorov-
1033 Smirnov test for performance in not motion corrected t-series vs. motion corrected t-series from LIV:
1034 $p = 0.54$ for F-1, $p = 0.15$ for Precision, $p = 0.38$ for Recall, $N = 13$ t-series. Results of Kolmogorov-
1035 Smirnov test for performance in not motion corrected t-series vs. motion corrected t-series from CA1
1036 jRCaMP1a: $p = 0.16$ for F-1, $p = 0.20$ for Precision, $p = 0.20$ for Recall, $N = 12$ t-series. Results of
1037 Kolmogorov-Smirnov test for performance in not motion corrected vs. motion corrected t-series from
1038 CA1 GCaMP6f: $p = 0.18$ for F-1, $p = 0.28$ for Precision, $p = 0.22$ for Recall, $N = 12$ t-series. **E)** F-1
1039 values as a function of the fraction of the total length of the t-series for not-motion corrected data (N
1040 = 13 t-series for LIV, $N = 24$ t-series for CA1, including $N = 12$ t-series for CA1 jRCaMP1a and N
1041 = 12 t-series for CA1 GCaMP6f t-series).



1043 **Figure 4. Fast extraction of fluorescence traces using CITE-On.** **A)** Median projection showing
1044 representative FOVs from the LIV GCaMP6s (left) and the CA1 jRCaMP1a (right) datasets. True
1045 positive bounding boxes for five CITE-On identified cells in each FOV are shown. **B)** Left: the five
1046 cells indicated in the LIV t-series displayed in (A) are shown at an expanded spatial scale. Right:
1047 corresponding background subtracted fluorescence traces. **C)** Same as in B but for the CA1 t-series
1048 in A. **D)** Lower-left triangle: cross correlation matrix for all functional traces extracted from true
1049 positive detection in the LIV GCaMP6s t-series displayed in (A). Upper-right triangle: corresponding
1050 dendrogram sorting. The left matrix shown signals before background subtraction. The right matrix
1051 after background subtraction. **E)** Same as in (D), but for the CA1 jRCaMP1a t-series shown in (A).

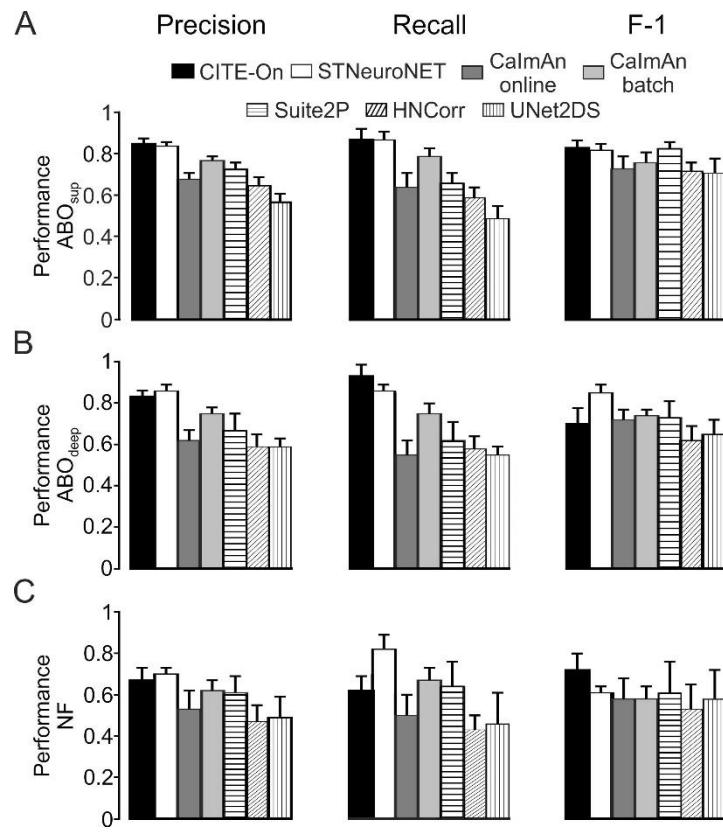


1052

1053 **Figure 5. CITE-On offline performance on never-before-seen data. A-E)** Best parameter search

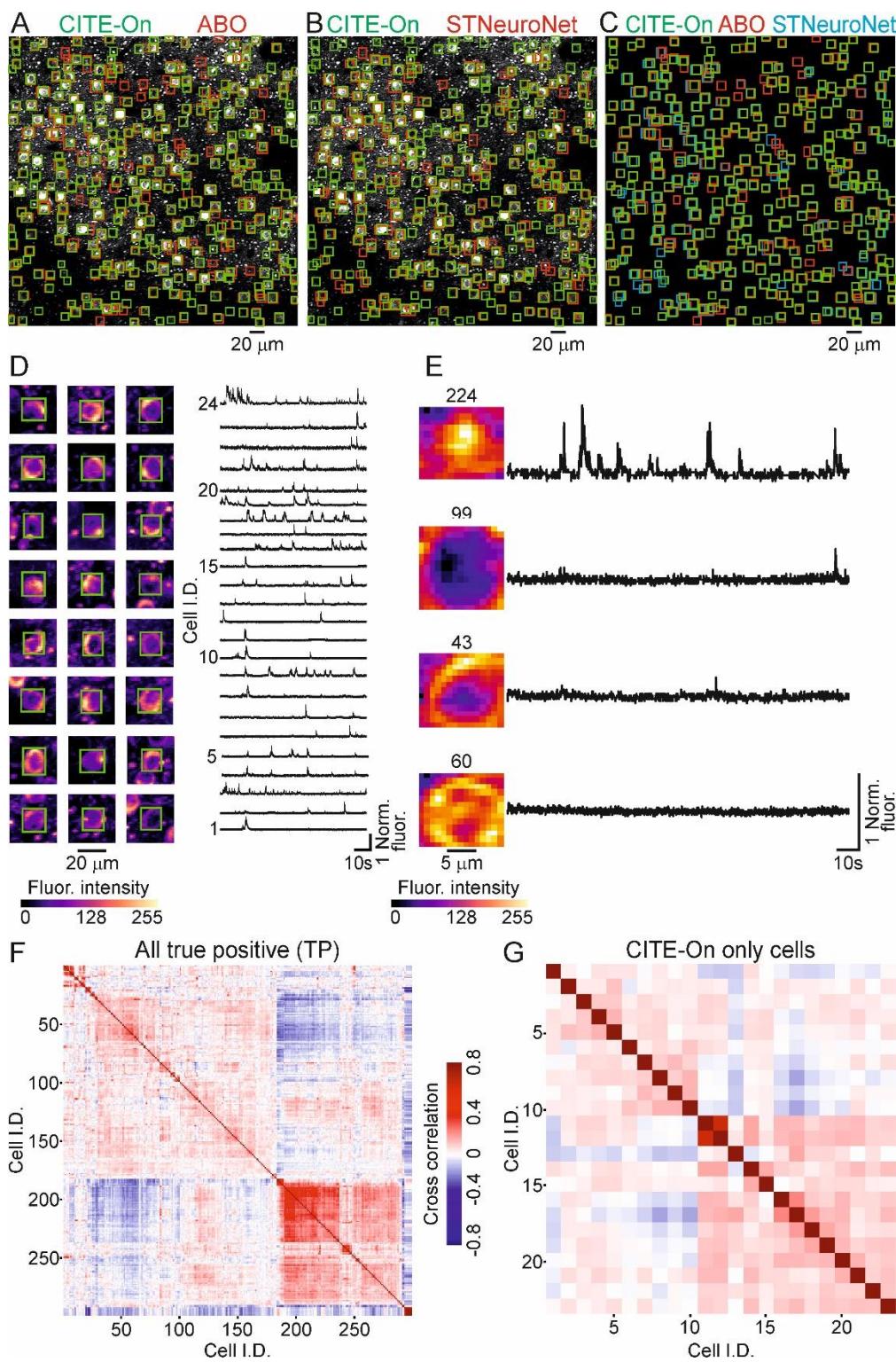
1054 for frame upscaling: F-1 score (pseudocolor) as a function of upscaling factor (vertical axis) and score

1055 threshold (horizontal axis) for the ABO_{sup} (A), ABO_{deep} (B), NF_{train} (C), NF_{test} (D) and VPM (E)
1056 datasets. The maximal F-1 is indicated with the black rectangle. The pseudocolor scale in (E) applies
1057 to (A-D). **F**) Performance of Precision (grey), Recall (white), and F-1 (black) for all t-series in the
1058 ABO_{sup} ($N = 9$), ABO_{deep} ($N = 10$), NF_{train} ($N = 19$), NF_{test} ($N = 9$), and VPM ($N = 9$) datasets.



1059

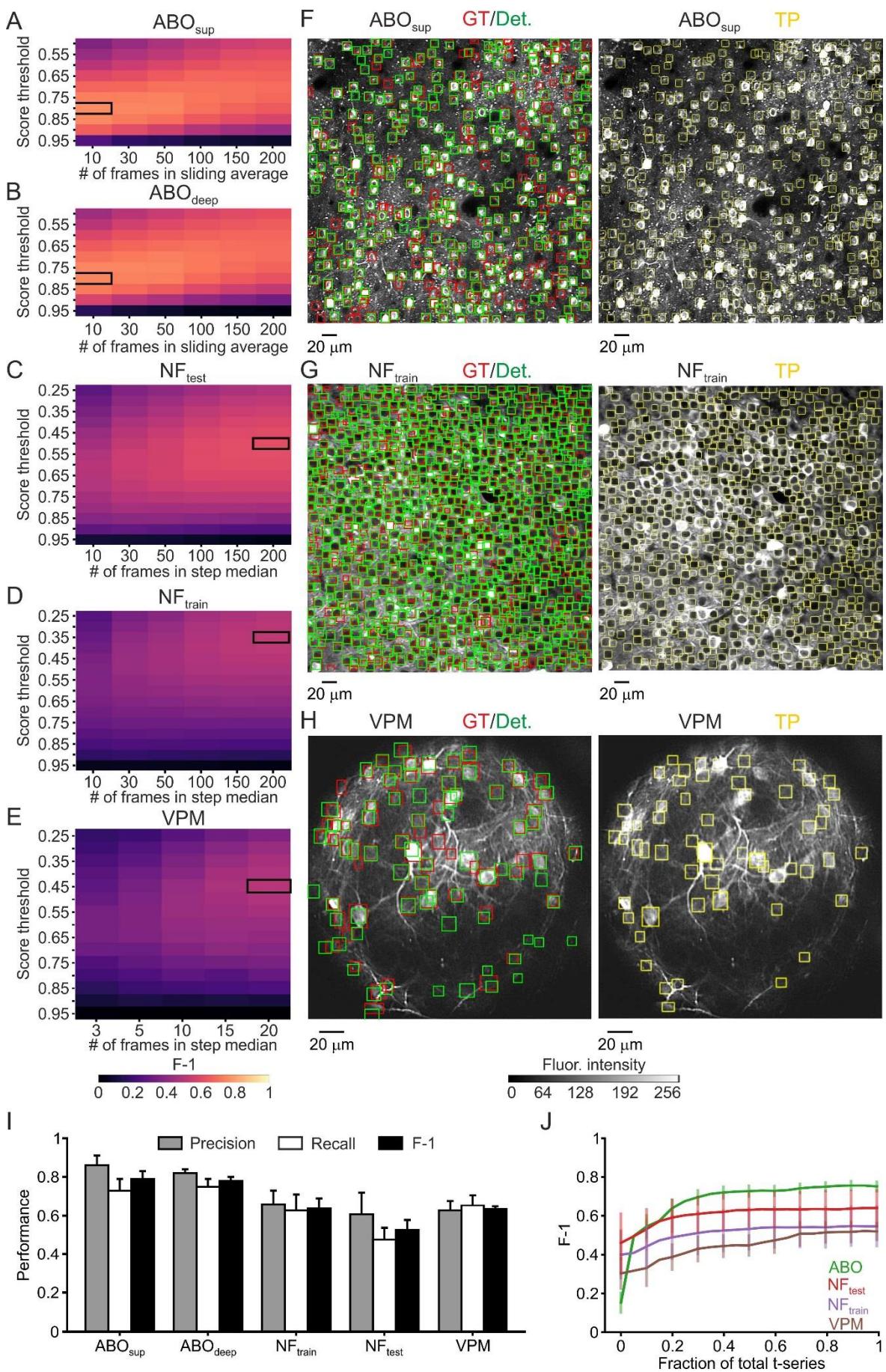
1060 **Figure 6. CITE-On offline performance compared with state-of-the-art methods.** A-C) Offline
 1061 performance of CITE-On (black) compared to STNeuroNET (white), CalmAn On Line (dark grey),
 1062 CalmAn Batch (light grey), Suite2P (horizontal line), HNCorr (tilted line) and UNet2DS (vertical
 1063 line). Precision (left), Recall (middle), and F-1 (right) are shown. Performance is evaluated on ABO_{sup}
 1064 (A, N = 9 t-series), ABO_{deep} (B, N = 10 t-series), and NF_{train} + NF_{test} (NF, C, N = 28 t-series).



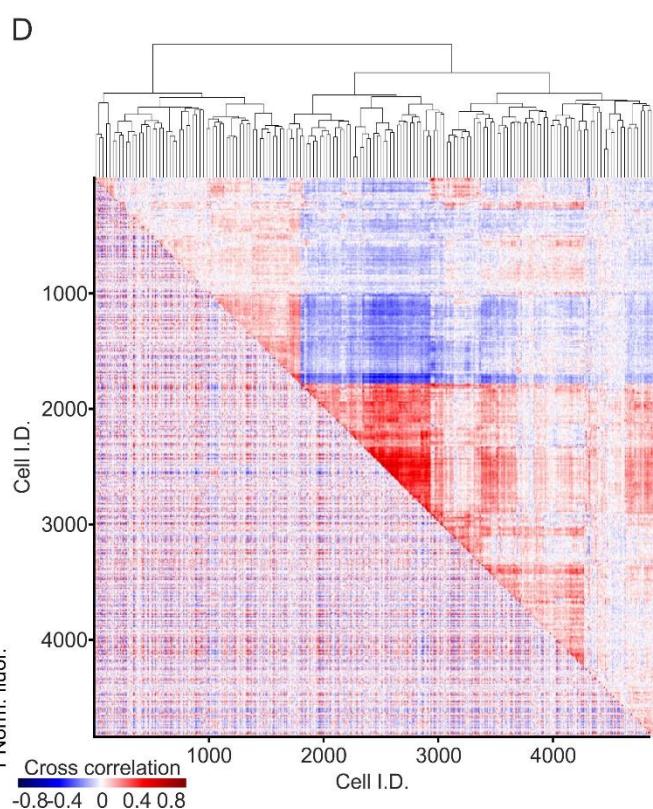
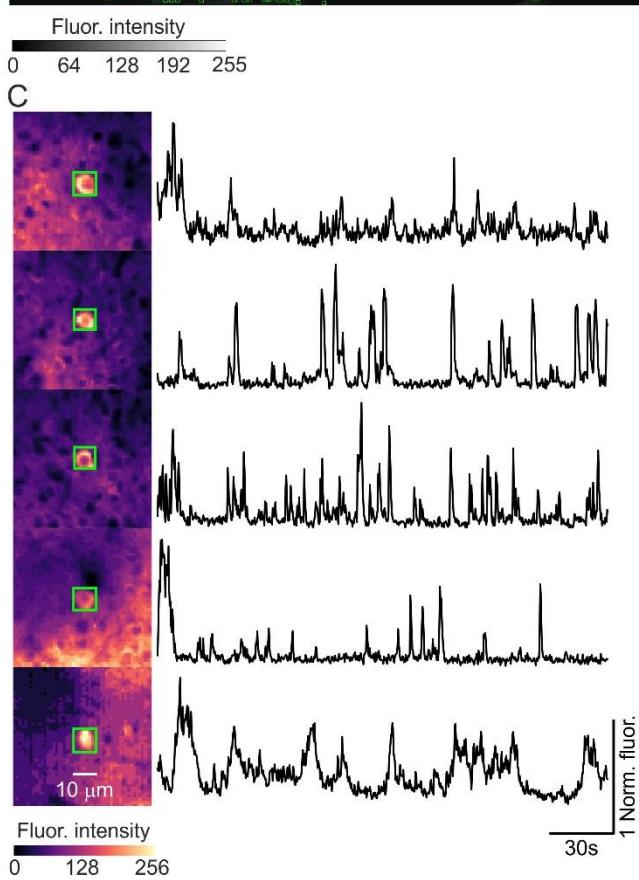
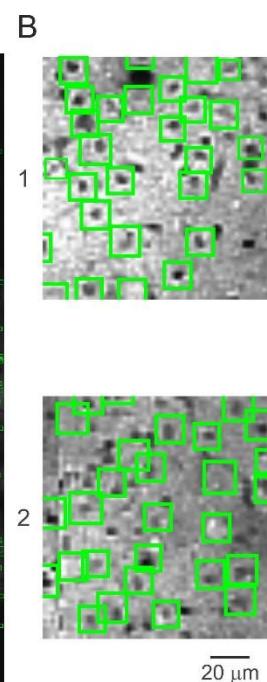
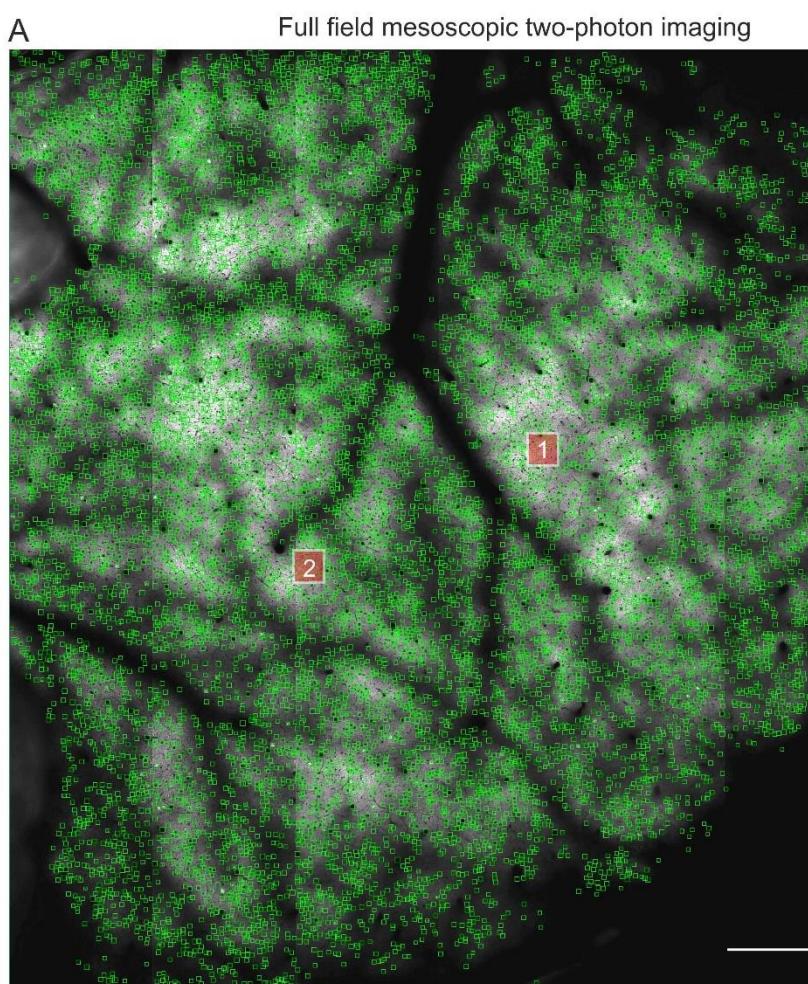
1065

1066 **Figure 7. CITE-On data processing of never-before-seen recordings.** A)
 1067 representative t-series from the ABO dataset showing GCaMP6f expressing cortical neurons. CITE-
 1068 On true positives (CITE-On, green) and true positives provided by the Allen Brain Observatory

1069 (ABO, red) are shown. **B**) Same as in (A) with CITE-On true positives (green) and STNeuroNET true
1070 positives (red). **C**) Superposition of CITE-On (green), ABO (red), and STNeuroNET (cyan) true
1071 positives. **D**) Left: 24 representative cells detected by CITE-On and identified as true positives in
1072 ABO and STNeuroNET. The CITE-On-identified bounding box is represented in green. Right:
1073 corresponding CITE-On-extracted fluorescence traces. **E**) Same as in (D) for four representative
1074 CITE-On-only cells. These fours cells were not counted in the GT of ABO and STNeuroNET GT,
1075 either as true or false positives. **F**) Lower-left triangle: cross correlation matrix for all functional traces
1076 extracted from true positive detection in the t-series displayed in (A). Upper-right triangle:
1077 corresponding dendrogram sorting. The pseudocolor scale indicates the cross correlation value. **G**)
1078 same as in (F), but for the CITE-On-only true positive cells. Pseudocolor scale as in (F).

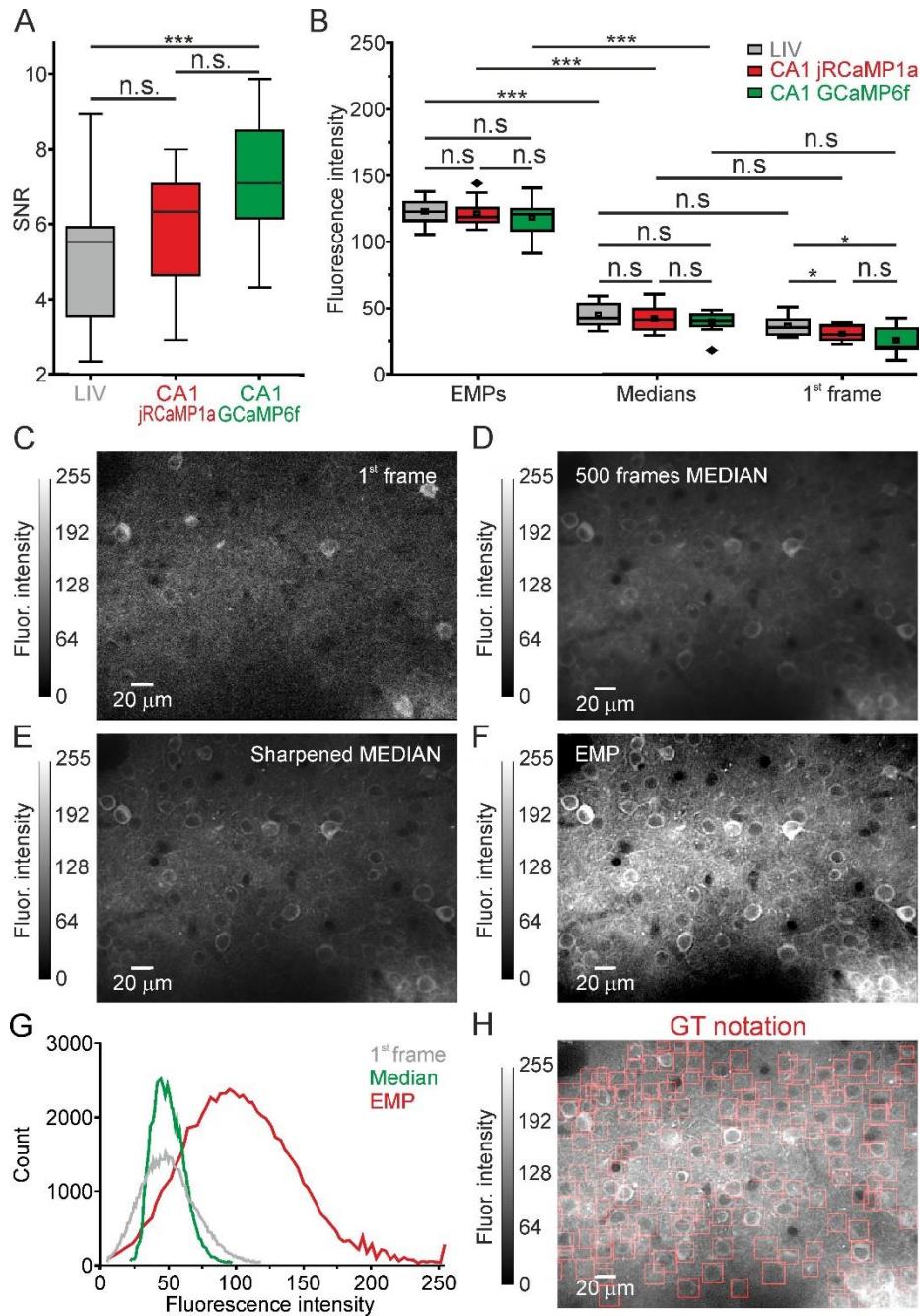


1080 **Figure 8. CITE-On online performance on never-before-seen datasets.** **A-E)** Best parameter
1081 search for frame downsampling: F-1 score (pseudocolor) as a function of score threshold (vertical
1082 axis) and number of frames (horizontal axis) for the ABO_{sup} (A), ABO_{deep} (B), NF_{train} (C), NF_{test} (D)
1083 and VPM (E) datasets. The maximal F-1 is indicated with the black rectangle. The pseudocolor scale
1084 in (E) applies to (A-D). For the ABO_{sup} , ABO_{deep} datasets the sliding average frame downsampling
1085 approach was used, while for the NF_{test} , NF_{train} , and VPM datasets the step average approach was
1086 implemented. **F-H)** Left: median projection of a representative t-series from the ABO_{sup} (F), NF_{train}
1087 (G), and VPM (H) datasets. GT (red) and online CITE-On detections (green bounding boxes) are
1088 shown. Right: bounding boxes (yellow) corresponding to true positives are shown. The greyscale in
1089 H applies also to (F-G). **I)** Online detection performance of Precision (grey), Recall (white), and F-1
1090 (black) for all t-series in the ABO_{sup} ($N = 9$), ABO_{deep} ($N = 10$), NF_{train} ($N = 19$), NF_{test} ($N = 9$), and
1091 VPM ($N = 9$) datasets. **J)** F-1 as a function of the fraction of processed t-series for ABO (green, $N =$
1092 19 t-series), NF_{test} (red, $N = 9$ t-series), NF_{train} (purple, $N = 19$ t-series), and VPM (brown, $N = 9$ t-
1093 series) datasets. 10 frames sliding averages for ABO; detection rate, 5 Hz. Step median of 20 frames
1094 and 200 frames for VPM and NF datasets; detection rate, 0.3 Hz and 0.035 Hz for SPM and NF
1095 datasets, respectively.



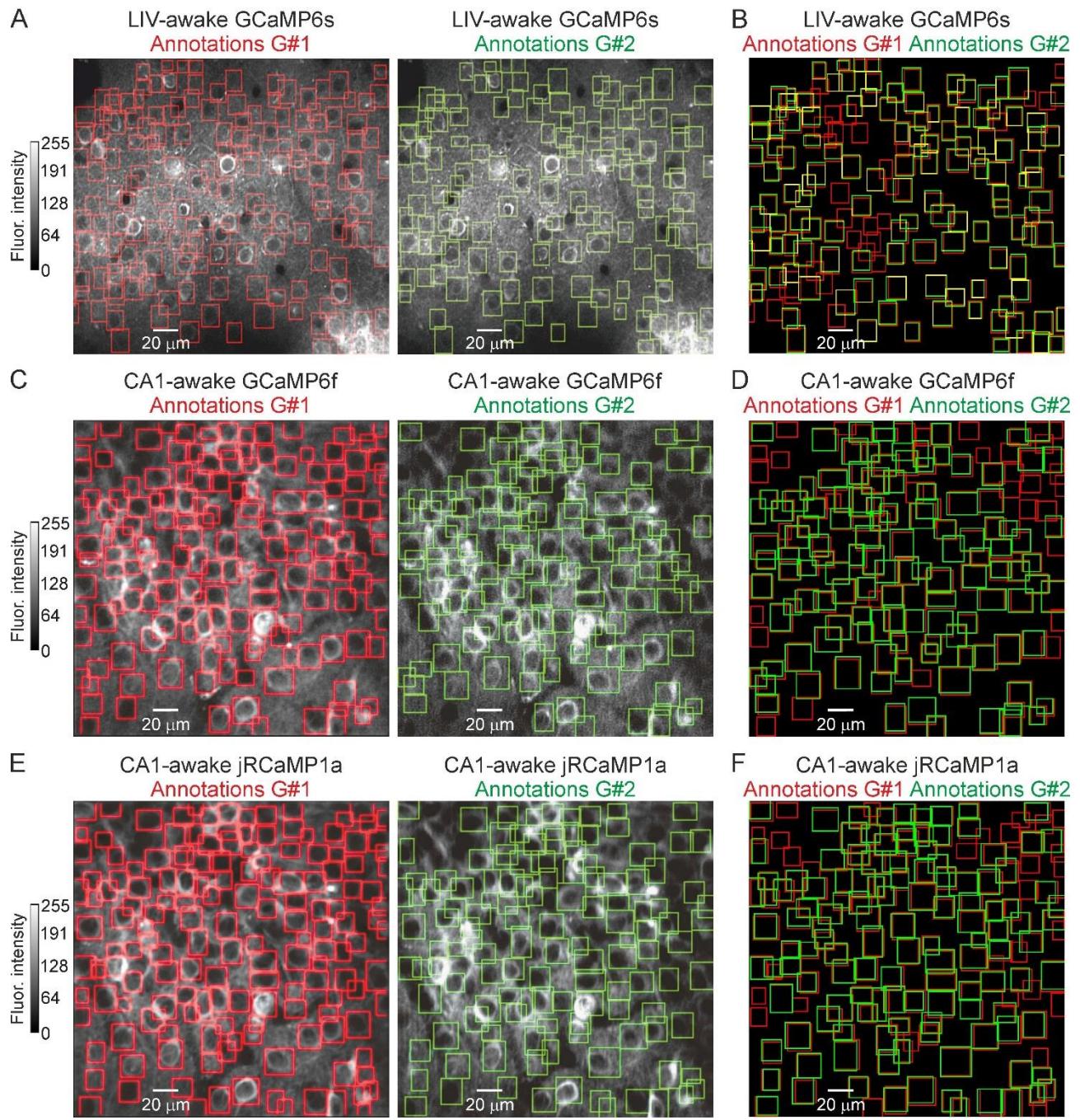
1097 **Figure 9. CITE-On analysis of mesoscopic two-photon imaging t-series.** **A-B)** Median projection
1098 of a mesoscopic imaging t-series showing GCaMP6s expressing neurons (mesoscopic data from ⁷).
1099 Green boxes indicate cells detected by CITE-On (total: 4,842 cells). Two regions are highlighted by
1100 yellow squares and shown at an expanded spatial scale in (B). Greyscale in (A) applies also to (B).
1101 **C)** Left: five representative cells detected by CITE-On. Right: corresponding CITE-On-extracted
1102 fluorescence traces in the first 230 s of the t-series. **D)** Cross correlation matrix (bottom-left triangle)
1103 calculated on the background subtracted traces extracted by CITE-On on all detected cells in the first
1104 7,000 frames and relative dendrogram (top-right triangle).

Supplementary material

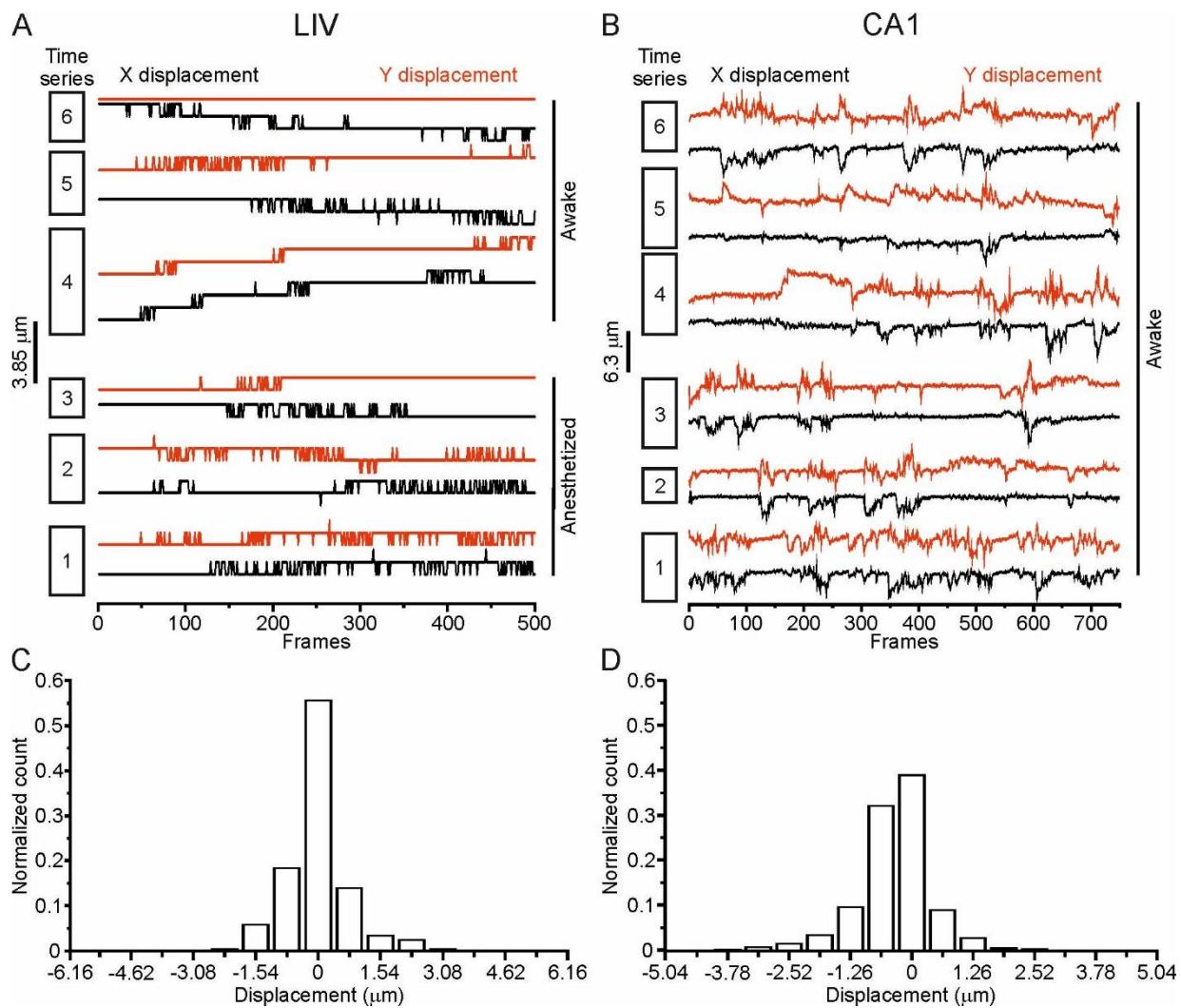


Supplementary Figure 1. Processing imaging t-series for GT annotation. **A)** SNR values across all dataset t-series (both training and validation). Unpaired Student's *t*-test, $p = 0.17$ between LIV vs. CA1 jRCaMP1a; $p = 0.009$ for LIV vs. CA1 GCaMP6f; $p = 0.26$ for CA1 GCaMP6f vs. CA1 jRCaMP1a. N = 121 for LIV, N = 38 for CA1 jRCaMP1a and for CA1 GCaMP6f. **B)** Fluorescence intensity for EMPs images, median projections of time series, and individual first frames of acquisitions for LIV (grey) and CA1 (GCaMP6f: green and jRCaMP1a: red) datasets. For EMPs, unpaired Student's *t*-test: $p = 0.77$ for LIV vs. CA1 jRCaMP1a; $p = 0.62$ for LIV vs. CA1 GCaMP6f; $p = 0.78$ for CA1 GCaMP6f vs. CA1 jRCaMP1a. For individual frames, Wilcoxon signed rank test: $p = 0.38$ for LIV vs. CA1 jRCaMP1a; $p = 0.56$ for LIV vs. CA1 GCaMP6f; $p = 0.68$ for CA1

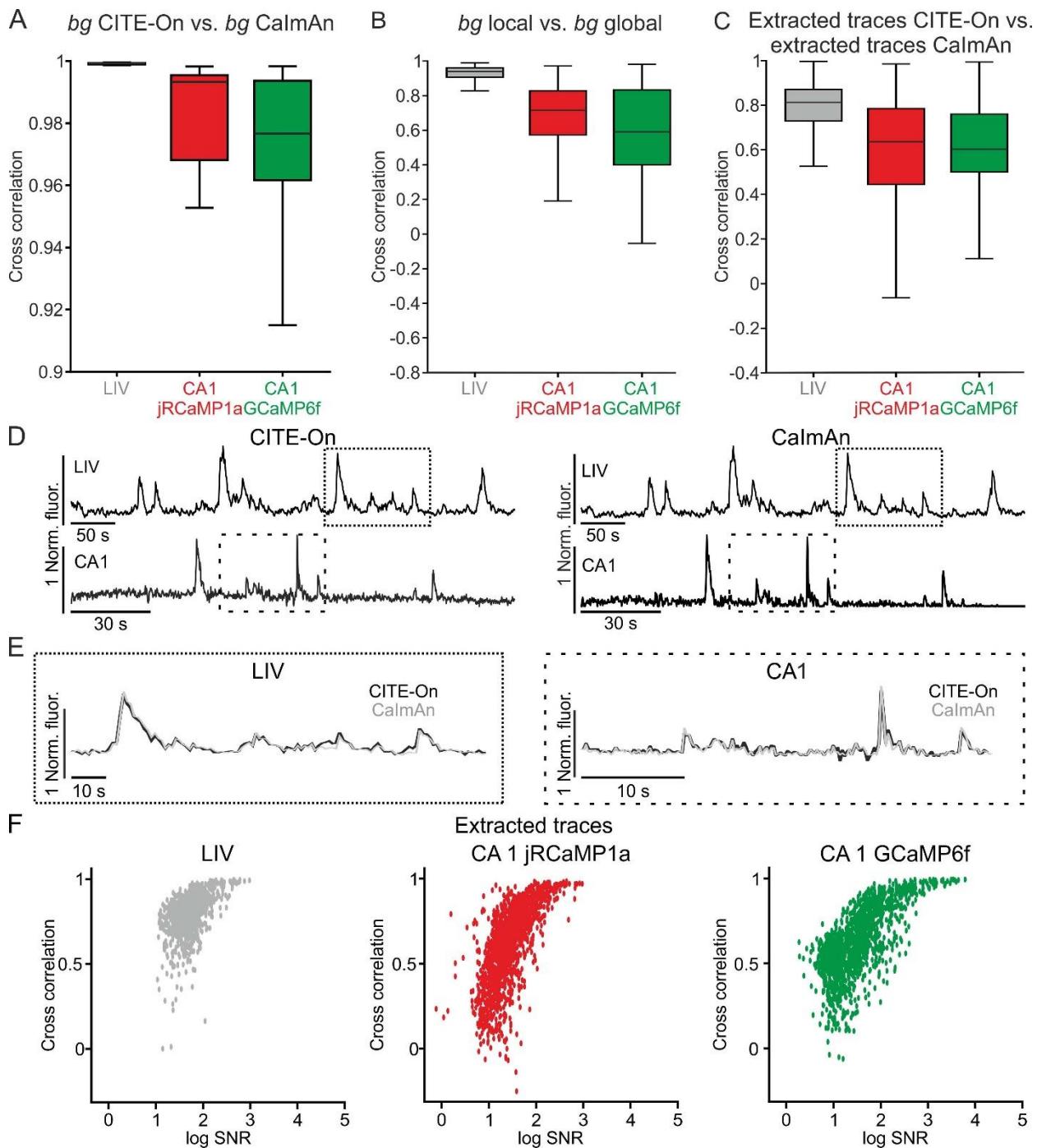
jRCaMP1a vs. CA1 GCaMP6f. For EMPs vs. median values, median values for LIV were significantly larger compared with individual frames of CA1 jRCAMP1a and CA1 GCaMP6f: Wilcoxon signed rank test, $p = 2.5\text{E-}5$ for LIV median vs. 1st frame CA1 jRCaMP1a; $p = 8.6\text{E-}6$ for LIV median vs. 1st frame CA1 GCaMP6f; $p = 0.082$ for LIV median vs. LIV 1st frame. Wilcoxon signed rank test $p = 3.5\text{E-}5$ for LIV EMP vs. LIV median; $p = 4.8\text{E-}8$ for LIV EMP vs. LIV 1st frame; $p = 5.6\text{E-}6$ for LIV EMP vs. CA1 jRCaMP1a median; $p = 4.8\text{E-}8$ for LIV EMP Vs CA1 jRCaMP1a 1st frame; $p = 5.6\text{E-}6$ for LIV EMP vs. CA1 GCaMP6f median; $p = 6.9\text{E-}8$ for LIV EMP vs. CA1 GCaMP6f 1st frame). N = 121 for LIV, N = 38 for CA1 jRCaMP1a and for CA1 GCaMP6f. **C-F**) Individual frame (C), median projection (D), global sharpened image (E), and EMP (F) for a representative LIV t-series. **G**) Distribution of absolute fluorescence intensity values for individual frames (grey), median projection (green), sharpened median projection (cyan) and EMP images (red) for a representative LIV t-series. **H**) GT annotation on a representative EMP.



Supplementary Figure 2. Graders' annotations on LIV and CA1 validation t-series. **A)** EMP image from a representative LIV t-series. Bounding boxes generated by grader #1 are shown in red in the left panel. Those generated by grader #2 are shown in green in the right panel. **B)** Superposition of the bounding boxes generated by grader #1 (red) and grader #2 (green). **C-D)** Same as in (A-B) for a representative CA1 GCaMP6f t-series. **E-F)** Same as in (A-B) for a representative CA1 jRCaMP1a t-series.

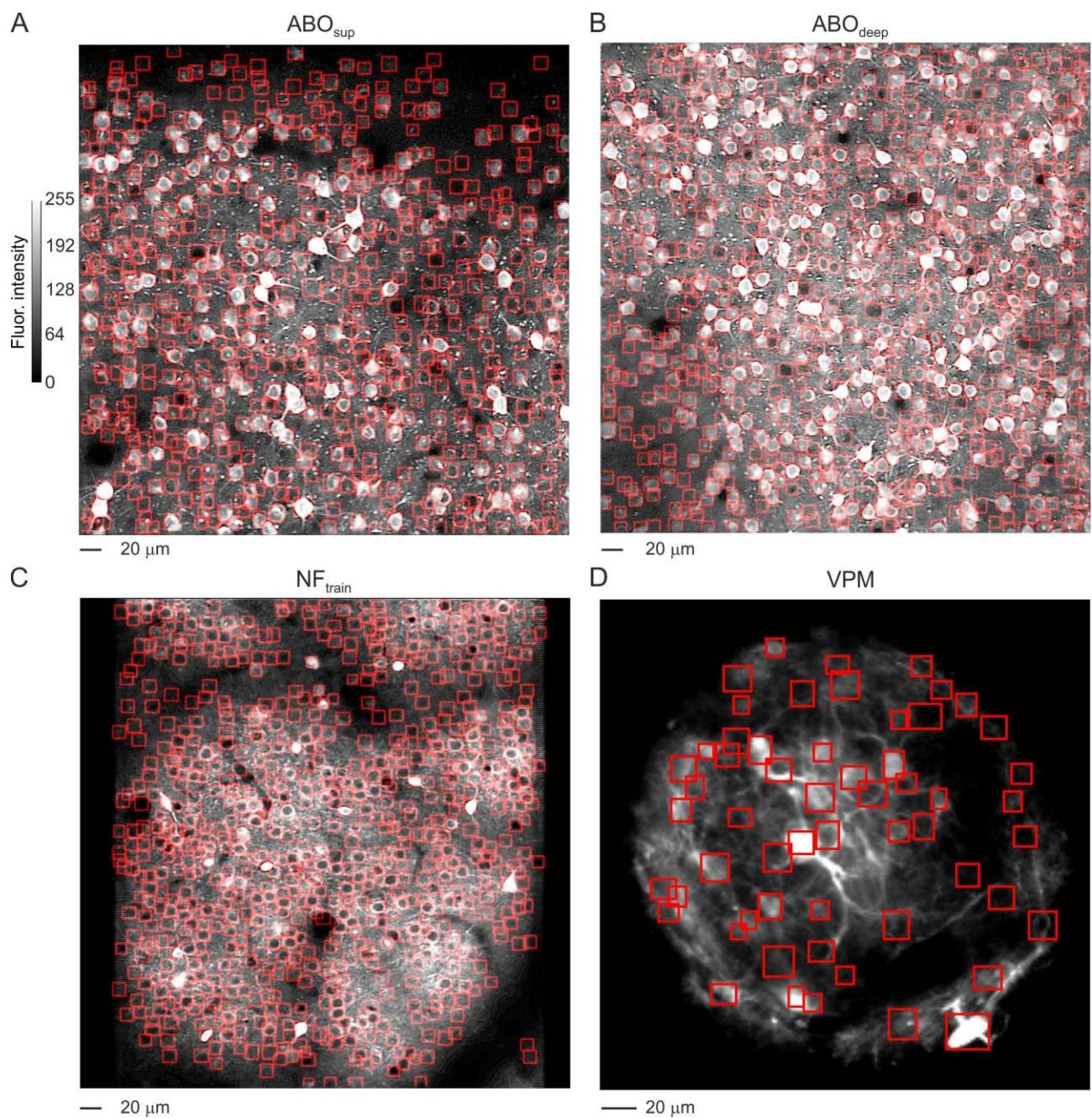


Supplementary Figure 3. Motion artefacts in the LIV and CA1 validation datasets. **A)** X, Y displacement of the FOV (black, X; red, Y) expressed in microns as observed across frames of six representative 500 frame-long LIV t-series (traces 1-3 from anesthetized animals, traces 4-6 from awake mice). **B)** Same as in (A) but for six representative 750 frame-long CA1 jRCaMP1a acquisitions. **C)** Percentage of total X, Y displacements in LIV validation t-series (N = 13). **D)** Same as in (C) for CA1 jRCaMP1a validation t-series (N = 12).

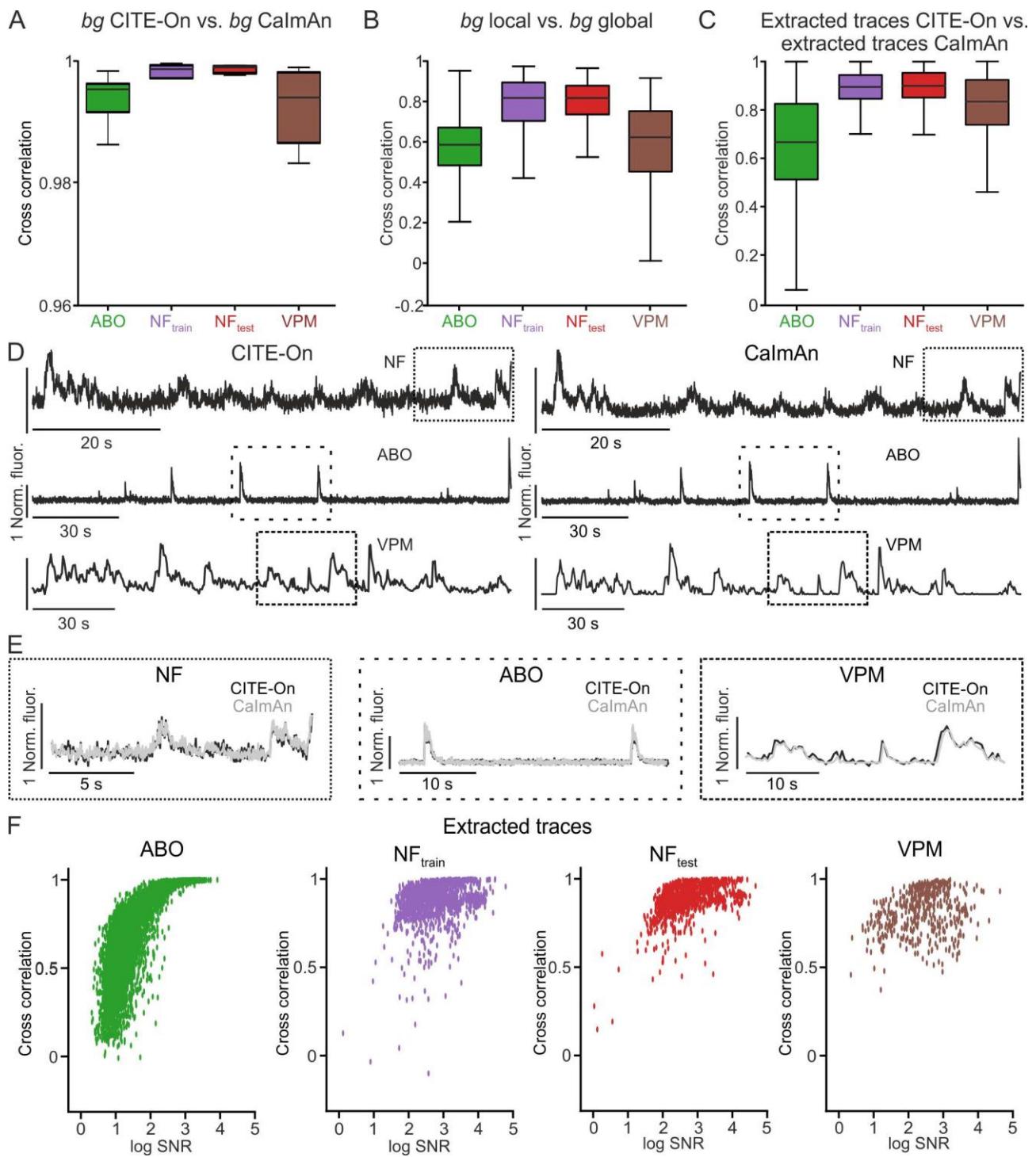


Supplementary Figure 4. Trace extraction from the validation dataset: CITE-On vs. seeded-CaImAn. **A)** Boxplots showing cross correlation values between the background obtained using CITE-On and seeded-CaImAn for the LIV, CA1 jRCaMP1a, and CA1 GCaMP6f datasets ($N = 13$, $N = 12$, and $N = 12$ t-series, respectively). **B)** Cross correlation values between global and local background signals computed with CITE-On. **C)** Cross correlation of background-subtracted functional traces extracted with CITE-On and with seeded-CaImAn for all true positive detected identities in the LIV, CA1 jRCaMP1a, and CA1 GCaMP6f datasets. **D)** Representative background subtracted functional traces extracted with CITE-On (left) and seeded-CaImAn (right) from LIV (top) and CA1 jRCaMP1a (bottom) t-series. **E)** Representative traces extracted with CITE-On (black) and CalmAn (grey) are shown superimposed for LIV (left) and CA1 (right). **F)** Cross correlation of

background-subtracted functional traces extracted with CITE-On and with seeded-CaImAn as a function of the cell's SNR for all true positive identities in the LIV (left), CA1 jRCaMP1a (middle), and CA1 GCaMP6f (right) datasets. Each dot represents a cell detected by CITE-On (see Supplementary Table 1).

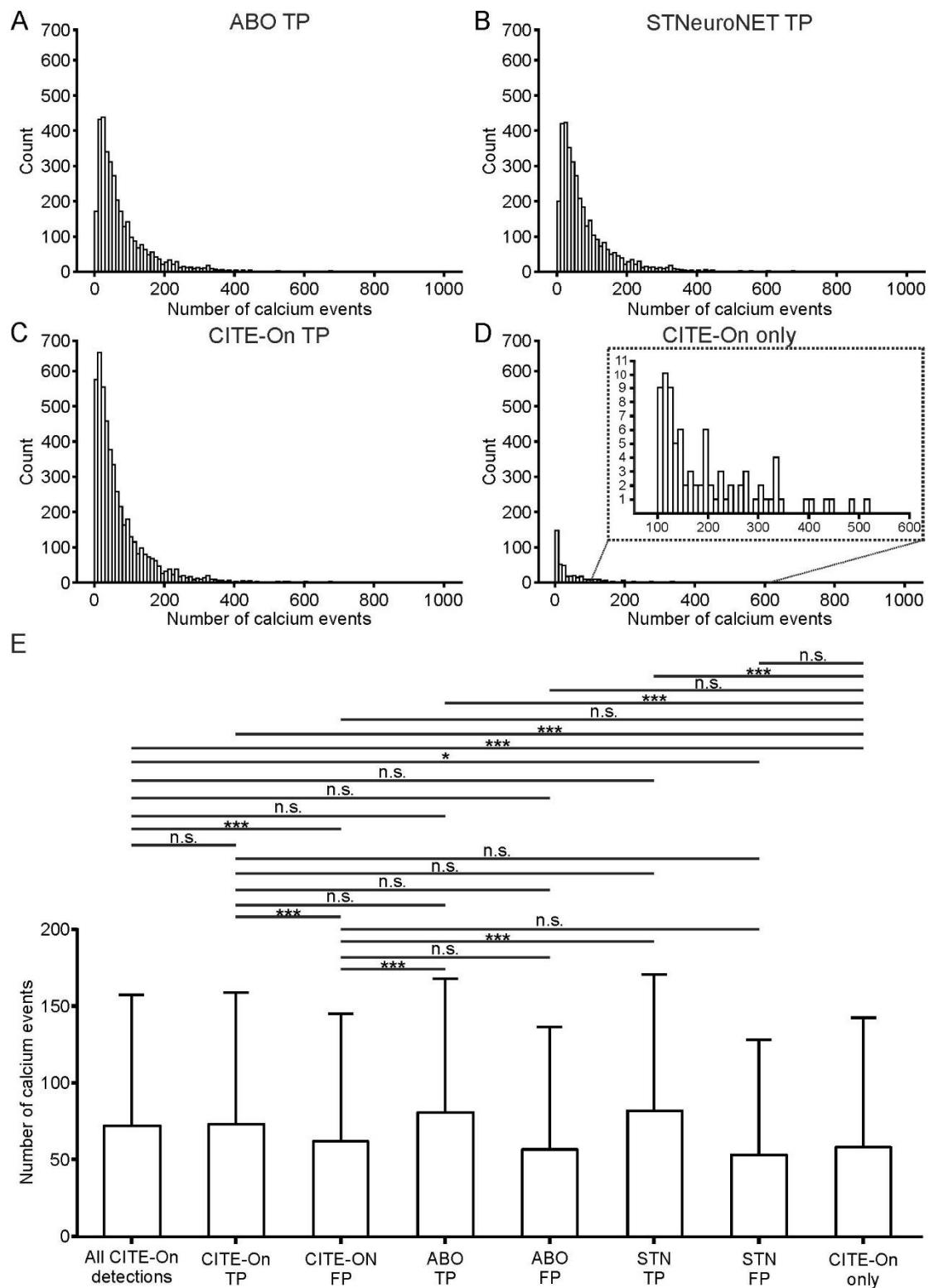


Supplementary Figure 5. GT bounding box annotation of publicly available datasets. A-D)
 Median projection of t-series from ABO_{sup} (A), ABO_{deep} (B), NF_{train} (C), and VPM (D) datasets. GT bounding boxes are shown in red.



Supplementary Figure 6. Trace extraction from publicly available datasets: CITE-On vs. seeded-CaImAn. **A)** Boxplots showing cross correlation values between the background obtained using CITE-On and seeded-CaImAn for the for ABO, NF_{train}, NF_{test}, and VPM data ($N = 19$, $N = 19$, $N = 9$, and $N = 9$ t-series, respectively). **B)** Cross correlation values between global and local background signals computed with CITE-On. **C)** Cross correlation of background-subtracted functional traces extracted with CITE-On and with seeded-CaImAn for the various datasets. **D)** Representative background subtracted functional traces extracted with CITE-On (left) and seeded-CaImAn (right) for NF (top), ABO (middle) and VPM (bottom) acquisitions. **E)** Representative traces

extracted with CITE-On (black) and CaImAn (grey) are shown superimposed for NF (left), ABO (middle), and VPM (right). **F**) Cross correlation of background-subtracted functional traces extracted with seeded-CaImAn and CITE-On as a function of the cell's SNR for all true positive identities in the ABO (leftmost), NF_{train} (middle left), NF_{test} (middle right) and VPM (rightmost) acquisitions. Each dot represents a cell detected by CITE-On (see Supplementary Table 1).



Supplementary Figure 7. Calcium activity in CITE-On, ABO, and STNeuroNET detections. **A-D**) Distribution of the number of calcium events in true positive detection in a representative ABO t-series analyzed by ABO (A), STNeuroNET (B), CITE-On (C). (D) shows the distribution of CITE-On-only true positive detections. Data refers to the corresponding GT annotation provided by the consensus GT for CITE-On, ABO, and STNeuroNET. **E)** Average (and s.d.) number of detected

calcium events *per* cell for all CITE-On detections, CITE-On true positives (CITE-On TP), CITE-On false positives (CITE-On FP), ABO true positives (ABO TP), ABO false positives (ABO FP), STNeuroNET true positives (STN TP), STNeuroNET false positives (STN FP), and CITE-On only detections. Data from N = 19 t-series of the ABO dataset. Results of Kolmogorov-Smirnov test: p = 0.00012 for All cells vs. CITE-On-only cells; p = 1.00 for All cells vs. CITE-On True Positives; p = 0.00022 for All cells vs. CITE-On False Positives; p = 0.99 for All cells vs. ABO True Positives; p = 0.078 for All cells vs. ABO False Positives; p = 0.99 for All cells vs. STN True Positives; p = 0.036 for All cells vs. STN False Positives; p = 0.00041 for CITE-On-only Cells vs. CITE-On True Positives; p = 0.99 for CITE-On-only Cells vs. CITE-On False Positives; p = 0.0012 for CITE-On-only Cells vs. ABO True Positives; p = 0.28 for CITE-On-only Cells vs. ABO False Positives; p = 0.00073 for CITE-On-only Cells vs. STN True Positives; p = 0.36 for CITE-On-only Cells vs. STN False Positives; p = 0.00073 for CITE-On True Positives vs. CITE-On False Positives; p = 0.99 for CITE-On True Positives vs. ABO True Positives; p = 0.11 for CITE-On True Positives vs. ABO False Positives; p = 0.99 for CITE-On True Positives vs. STN True Positives; p = 0.078 for CITE-On True Positives vs. STN False Positives; p = 0.0022 for CITE-On False Positives vs. ABO True Positives; p = 0.46 for CITE-On False Positives vs. ABO False Positives; p = 0.0012 for CITE-On False Positives vs. STN True Positives; p = 0.70 for CITE-On False Positives vs. STN False Positives; p = 0.28 for ABO True Positives vs. ABO False Positives; p = 1.00 for ABO True Positives vs. STN True Positives; p = 0.15 for ABO True Positives vs. STN False Positives; p = 0.21 for ABO False Positives vs. STN True Positives; p = 0.99 for ABO False Positives vs. STN False Positives; p = 0.11 for STN True Positives vs. STN False Positives. Total number of detected cells: N = 5482, for All CITE-On detections; N = 4934, for CITE-On TP; N = 548, for CITE-On FP; N = 3516, for ABO TP; N = 1966, for ABO FP; N = 3606, for STN TP; N = 1876, for STN FP; N = 439, for CITE-On only.

Supplementary tables

Dataset	# of t-series	# of detections grader 1	# of detections grader 2	# of detections overlap	# of detections grader 1 only	# of detections grader 2 only	# of ground truth
ABO	19	6225	5456	5443	782	13	6238
NF _{train}	19	11171	9772	9630	1541	142	11313
NF _{test}	9	3187	2814	2793	394	21	3208
VPM	9	517	470	444	73	26	543
LIV _{test}	13	979	868	861	118	7	986
CA1 _{test} jRCaMP1a	12	884	809	795	89	14	898
CA1 _{test} GCaMP6s	12	1000	927	916	84	11	1011
LIV _{train}	118	7221	6489	6356	865	133	7354
CA1 _{train} jRCaMP1a	21	2123	1856	1833	290	23	2146
CA1 _{train} GCaMP6s	21	2218	1963	1940	278	23	2241

Supplementary table 1. Dataset annotation and ground truth generation. Two graders manually annotated the LIV, CA1, ABO, NF, and VPM datasets. The table reports the number of t-series, the

number of detections by grader 1 and 2, the number of overlapping detections, the number of detections exclusively produced by grader 1 and grader 2, and the number of detections in the consensus ground truth.

Dataset	mAP	s.d	N	F-1	s.d	N	Precision	s.d	N	Recall	s.d	N
ABO	0.80	0.10	19	0.93	0.02	19	0.998	0.004	19	0.87	0.04	19
NF _{train}	0.71	0.14	19	0.92	0.03	19	0.99	0.02	19	0.86	0.05	19
NF _{test}	0.72	0.10	9	0.93	0.02	9	0.99	0.01	9	0.88	0.04	9
VPM	0.64	0.13	9	0.90	0.03	9	0.95	0.03	9	0.86	0.04	9
LIV _{test}	0.81	0.11	13	0.93	0.03	13	0.99	0.01	13	0.88	0.05	13
CA1 _{test} jRCaMP1a	0.76	0.07	12	0.94	0.03	12	0.98	0.02	12	0.90	0.04	12
CA1 _{test} GCaMP6s	0.82	0.08	12	0.95	0.01	12	0.99	0.01	12	0.92	0.03	12
LIV _{train}	0.73	0.09	118	0.93	0.02	118	0.98	0.02	118	0.88	0.05	118
CA1 _{train} jRCaMP1a	0.74	0.08	21	0.92	0.02	21	0.99	0.01	21	0.87	0.04	21
CA1 _{train} GCaMP6s	0.76	0.07	21	0.93	0.02	21	0.99	0.01	21	0.87	0.04	21

Supplementary table 2. Evaluation of consensus GT. The table reports the average (and relative s.d.) of the mean average precision (mAP), the F-1 score, the precision and recall across the t-series of the various datasets annotated by the different graders.

Offline												
Dataset	F-1	s.d	N	Precision	s.d	N	Recall	s.d	N	# of TP	# of FP	# of FN
ABO	0.84	0.03	19	0.90	0.06	19	0.80	0.07	19	4934	548	1303
NF _{train}	0.63	0.1	19	0.63	0.09	19	0.62	0.12	19	6849	3932	4120
NF _{test}	0.74	0.08	9	0.70	0.07	9	0.51	0.06	9	2286	716	911
VPM	0.75	0.02	9	0.74	0.06	9	0.77	0.06	9	398	140	123
LIV _{test}	0.86	0.04	13	0.85	0.05	13	0.88	0.05	13	1486	257	205
CA1 _{test} jRCaMP1a	0.92	0.04	12	0.88	0.06	12	0.97	0.02	12	1618	222	55
CA1 _{test} GCaMP6s	0.93	0.02	12	0.91	0.03	12	0.97	0.02	12	1638	170	56

Supplementary table 3. CITE-On offline performance. Average (and relative s.d.) F-1, Precision, Recall, and number of true positives (TP), false positives (FP), and false negatives (FN) for LIV, CA1, ABO, NF and VPM datasets generated by the CITE-On offline pipeline.

Online

Dataset	F-1	s.d	N	Precision	s.d	N	Recall	s.d	N	# of TP	# of FP	# of FN
ABO	0.77	0.03	19	0.83	0.04	19	0.83	0.04	19	4573	978	1664
NF _{train}	0.55	0.11	19	0.56	0.11	19	0.56	0.11	19	6164	4906	5026
NF _{test}	0.64	0.10	9	0.68	0.12	9	0.68	0.12	9	2096	921	1101
VPM	0.52	0.50	9	0.60	0.11	9	0.60	0.11	9	244	181	277
LIV _{test}	0.69	0.04	13	0.72	0.05	13	0.72	0.05	13	1127	437	564
CA1 _{test} jRCaMP1a	0.73	0.07	12	0.74	0.07	12	0.74	0.07	12	1233	439	440
CA1 _{test} GCaMP6s	0.67	0.12	12	0.69	0.12	12	0.69	0.12	12	1145	493	549

Supplementary table 4. CITE-On online performance. Average (and relative s.d.) F-1, Precision, Recall, and number of true positives (TP), false positives (FP), and false negatives (FN) for LIV, CA1, ABO, NF and VPM datasets generated by the CITE-On online pipeline.

ABO Dataset	# of detections	s.d	N
All CITE-On	289	21	19
CITE-On TP	260	28	19
CITE-On FP	29	16	19
ABO TP	185	36	19
ABO FP	103	31	19
STNeuroNET TP	190	33	19
STNeuroNET FP	99	24	19
CITE-On only	23	15	19

Supplementary table 5. Detections in the ABO dataset. Average (and relative s.d.) number of detections in the ABO dataset reported by CITE-On, ABO, and STNeuroNET. TP, true positives; FP, false positives. Data relative to CITE-On were calculated online as the number of identities obtained at the end of the processing of each t-series then averaged across all t-series.

Supplementary movie legends

Supplementary movie 1. Online tracking of detected identities across the t-series. A representative t-series (# of frames, 750; frame rate, 3 Hz; no motion correction) from the CA1 validation dataset was processed online for cell detection (detection update rate, 10 Hz). Bounding boxes (colored squares) for active detections (i.e. neurons identified by CITE-On in the current frame) and past detections (i.e. neurons identified by CITE-On in previous frames) are represented with or without a central dot, respectively. Each identity was associated with a bounding box color which was retained across the t-series. The position and shape of each bounding box was updated in each frame according to the procedure described in the Results section.

Supplementary movie 2. Frame-by-frame manual annotation. Representative t-series (# of frames, 500; frame rate, 1.5 Hz) from the LIV training dataset (awake head-fixed *Scnn1a-cre* mouse expressing GCaMP6s) used for manual annotation of individual frames. The t-series was motion corrected. Brightness and contrast of each displayed frame of the t-series could be adjusted. Bounding boxes were manually positioned around each visible cell in each frame. Only a minority of neurons were identified in the frame-by-frame annotation process.

Supplementary movie 3. Manual annotation of EMP images. EMP image corresponding to the same t-series as in Supplementary movie 2. Brightness and contrast of the EMP image could be regulated to better visualize bright and dim cells. Bounding boxes (green) were manually defined around identified neurons.

Supplementary movie 4. Online dynamic segmentation and functional trace extraction. Top left: raw fluorescence of a representative cell from a t-series in the ABO dataset. The displayed cell was detected online using CITE-On. Top right: binary mask generated by CITE-On on each frame of the t-series displayed on the left. The white pixels represent the segmented pixels, black pixels are discarded. Bottom: functional fluorescence trace extracted online for the cell displayed on top.